Archetypal Analysis for Sparse Nonnegative Matrix Factorization: Robustness Under Misspecification

Kayhan Behdin*1 and Rahul Mazumder $^{\dagger 1,2}$

¹Operations Research Center, Massachusetts Institute of Technology ²Sloan School of Management, Massachusetts Institute of Technology

Abstract

We consider the problem of sparse nonnegative matrix factorization (NMF) with archetypal regularization. The goal is to represent a collection of data points as nonnegative linear combinations of a few nonnegative sparse factors with appealing geometric properties, arising from the use of archetypal regularization. We generalize the notion of robustness studied in Javadi and Montanari (2019) (without sparsity) to the notions of (a) strong robustness that implies each estimated archetype is close to the underlying archetypes and (b) weak robustness that implies there exists at least one recovered archetype that is close to the underlying archetypes. Our theoretical results on robustness guarantees hold under minimal assumptions on the underlying data, and applies to settings where the underlying archetypes need not be sparse. We propose new algorithms for our optimization problem; and present numerical experiments on synthetic and real datasets that shed further insights into our proposed framework and theoretical developments.

1 Introduction

Nonnegative Matrix Factorization (NMF) (Lee and Seung, 1999) is a well-known dimensionality reduction method where we represent a collection of data points as nonnegative linear combinations of a few nonnegative latent factors. Nonnegative factors are desirable from an interpretability standpoint in applications such as computational biology (Brunet et al., 2004; Kotliar et al., 2019), image processing (Kalayeh et al., 2014; Liu et al., 2011), text mining (Berry and Browne, 2005), and chemometrics (Lawton and Sylvestre, 1971), among others. Mathematically, given a $m \times n$ data matrix with nonnegative entries $\boldsymbol{X} \in \mathbb{R}^{m \times n}_{\geq 0}$ (the rows of \boldsymbol{X} correspond to the m samples and the columns the dimensions), NMF computes nonnegative lower-dimensional latent factors $\boldsymbol{W} \in \mathbb{R}^{m \times k}_{\geq 0}$. Here, k denotes the number of latent factors with k < m, n and we desire the factors to lead to a good approximation of the underlying data matrix: $\boldsymbol{X} \approx \boldsymbol{W}\boldsymbol{H}$, where

^{*}behdink@mit.edu

 $^{^\}dagger$ rahulmaz@mit.edu

rows of \mathbf{H} are representatives of the data and rows of \mathbf{W} denote the coefficient weights. In the simplest form, NMF can be formulated (Lin, 2007) as the following nonconvex optimization problem

$$\min_{\boldsymbol{H}.\boldsymbol{W}} \|\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H}\|_F^2 \quad \text{s.t.} \quad \boldsymbol{H} \in \mathbb{R}_{\geq 0}^{k \times n}, \boldsymbol{W} \in \mathbb{R}_{\geq 0}^{m \times k}$$
(1)

where, $\|\cdot\|_F$ denotes the Frobenius norm of a matrix.

Archetypal Analysis. The NMF problem is inherently under-determined due to scaling issues; and additional constraints should be imposed to make it well-defined. Archetypal Analysis (AA) due to Cutler and Breiman (1994) is a regularized variant of NMF where representatives of the data (aka the archetypes), given by the rows of \mathbf{H} , have an appealing geometric interpretation. The archetypes are chosen so that they belong to the convex hull of the data and the data is contained within their convex hull. In practice however, it may not be possible to find such archetypes (Javadi and Montanari, 2019). To this end, Javadi and Montanari (2019) propose a regularization scheme based on AA: Among all possible sets of candidate archetypes that represent the data with an acceptable accuracy, we select the one that is the closest to the convex hull of the data. Figure 1 (a), (b) illustrate the exact AA of Cutler and Breiman (1994) and regularized AA with a toy example 1 .

Robustness. If X admits an exact factorization of the form X = WH for some $W \in \mathbb{R}_{>0}^{m \times k}$ and $H \in \mathbb{R}_{>0}^{k \times n}$, Donoho and Stodden (2004) show that under the so-called separability assumption², it is possible to recover W and H from X (up to permutation and scaling of rows of H/columns of W). The separability assumption has been consequently generalized to less restrictive cases, including noisy settings. In particular, Arora et al. (2012) consider an approximately separable model, where the data is assumed to be a noisy version of a separable dataset. They show that in this model, under additional regularity conditions, a polynomialtime algorithm exists that finds a factorization that is close to the factorization of the noiseless data in a suitable metric. Their results are further improved by Recht et al. (2012), showing that noisy separable NMF can be solved by linear programs. Ge and Zou (2015) show that a relaxed version of separability, the so-called subset separability condition, suffices to achieve a good factorization from noisy data. Note that the separability/subset separability assumptions are usually hard to verify on real data; and in our development we do not make use of this assumption. Javadi and Montanari (2019) show that AA enjoys robustness to perturbation: under certain assumptions, the resulting solution from AA on the perturbed data is close to the underlying model in a suitable metric (as discussed in Section 2.2). In particular, this implies that at least one of the recovered archetypes is close to the underlying archetypes that contain and represent the noiseless data. In this paper, we generalize this notion to a stronger version of robustness which implies that each recovered archetype is close to some true archetype.

Some earlier works have considered NMF formulations robust to outliers (Chen et al., 2014; Kong et al., 2011). In this paper, we consider a different notion of robustness—unrelated to outliers arising in the context of robust statistics (Huber, 2004).

Sparsity. Generally speaking, due to nonnegativity constraints on the latent factors, NMF is known to produce sparse solutions, that is, \boldsymbol{H} and \boldsymbol{W} have some zero entries (Yang and Oja, 2010). A sparse representation of the data aids in interpretability and requires less storage space. This property of NMF has

¹See Appendix B for details of the example.

²The factorization X = WH is called separable if rows of X are among rows of H.

been utilized in different applications such as image processing (Hoyer, 2004), computational biology (Kim and Park, 2007), medical imaging (Woo et al., 2018), document clustering (Kim and Park, 2008) and audio processing (Virtanen, 2007). Several papers have proposed formulations of NMF with additional penalties and/or constraints to encourage enhanced sparsity in NMF—we refer to these as sparse NMF methods. Specifically, Hoyer (2004) consider a sparse NMF problem where they use a combination of the ℓ_1 and ℓ_2 penalty on the entries of \mathbf{H} . Penarz and Pernkopf (2012) add a constraint of the form $\|\mathbf{H}\|_0 \leq \ell$ to problem (1), where $\|\boldsymbol{H}\|_0$ is the ℓ_0 -pseudonorm, the number of nonzero entries of \boldsymbol{H} and ℓ is the desired sparsity level. Kim and Park (2007, 2008) add an ℓ_1 norm penalty on the entries of \mathbf{H} to the cost function of (1) to impose sparsity on H. To the best of our knowledge, there is limited theoretical work on sparsity in NMF—in particular, towards understanding the effect of sparsity constraints on the robustness of the representation returned by NMF, an aspect we study here. In this paper, we present a simultaneous analysis of sparsity and archetypal regularization in the form of Sparse AA (SAA). We study regularized AA (Javadi and Montanari, 2019) in the presence of additional sparsity constraints on H. In other words, we look for H such that it has a few nonzero entries (i.e., $\|\boldsymbol{H}\|_0$ is small), its rows describe the data well and are close to the convex hull of data. See Figure 1 (c) for a numerical illustration of this problem. In particular, we show that sparsity constrained AA leads to robust solutions, both in the weak sense of Javadi and Montanari (2019) and the stronger notion of robustness proposed here. An important feature of our analysis is that we do not assume the underlying archetypes are sparse—i.e., we can handle model misspecification—this makes our proofs different from existing work. We also discuss how noise and sparsity affect the robustness properties of the model.

Algorithms. Due to the bilinearity of the mapping $(W, H) \mapsto WH$, most formulations of NMF end up in a nonconvex optimization problem, although some convex formulations exist (Bach et al., 2008) when the dimension of the latent factors grows to infinity. Some basic approaches to these nonconvex problems include projected gradient methods (Lin, 2007), multiplicative update rules (Gonzalez and Zhang, 2005; Lee and Seung, 1999) and alternating optimization (Chu et al., 2004; Paatero and Tapper, 1994). More sophisticated algorithms for NMF have been proposed in recent years, for example see Gillis and Vavasis (2014); Leplat et al. (2019); Mizutani (2014). In this paper we present algorithms to obtain good solutions for the regularized AA problem with sparsity constraints. To this end, we present proximal block coordinate methods, and establish that they lead to a stationary point. We discuss a useful initialization scheme based on Mixed Integer Programming (MIP) (Bertsimas et al., 2016; Wolsey and Nemhauser, 1999) that leads to high-quality solutions. To further improve the quality of solutions available from our block coordinate procedure, we present local search based methods—to this end, our framework draws inspiration from the work of Beck and Eldar (2013); Hazimeh and Mazumder (2020) and adapts it to the setting of matrix factorization problems. Note that Abrol and Sharma (2020); Elhamifar et al. (2012); Mørup and Hansen (2012) have proposed algorithms for the original AA problem without sparsity constraints. In addition, prior work has extended the notion of separability arising from NMF to address the AA problem: For example, Damle and Sun (2017) use geometric interpretations of AA and separability to develop a new algorithm for NMF. Our numerical experiments on synthetic, and real datasets validate our theoretical results, and suggest the superiority of SAA over other popular sparse NMF methods.

Our Contributions. Our contribution in this paper can be summarized as follows:

• We generalize the robustness framework of Javadi and Montanari (2019) to notions of weak and strong

robustness, introduced in this paper—these notions differ in how they describe the proximity of our estimators to the underlying archetypes. Furthermore, we prove robust solutions are good representatives of the noiseless data.

- We show how sparsity and AA can be used together to produce sparse factors that are robust to noise and perturbation in the data. Our results apply to the mis-specified setting—i.e., situations where the underlying archetypes are not necessarily sparse.
- We present algorithms³ based on block proximal descent and local search, discuss MIP-based initialization strategies; and present convergence properties of our proposed algorithmic framework. We demonstrate via numerical experiments on synthetic and real-datasets the usefulness of our proposed approach.

Notation. For a matrix X, we let $X_{i,j}$, $X_{i,.}$ and $X_{.,j}$ denote the (i,j)-th element, i-th row and j-th column of X, respectively. With a slight abuse of notation, for indexed variables, we move the subscript to superscript: that is, the i-th row of X_0 is shown as $X_{i,.}^0$. The number of rows of a matrix X is denoted as #row(X). The convex hull of rows of the matrix $X \in \mathbb{R}^{m \times n}$ is denoted by

$$Conv(\boldsymbol{X}) = \left\{ \sum_{i=1}^{m} \alpha_i \boldsymbol{X}_{i,.} : \alpha_i \ge 0, \sum_{i=1}^{m} \alpha_i = 1 \right\}.$$

We define (i) the distance between a vector \boldsymbol{x} and the convex hull $\operatorname{Conv}(\boldsymbol{X})$ as $D(\boldsymbol{x}, \boldsymbol{X}) = \min_{\boldsymbol{v} \in \operatorname{Conv}(\boldsymbol{X})} \|\boldsymbol{x} - \boldsymbol{v}\|_2^2$ (ii) the distance between a set of points (i.e., the rows of \boldsymbol{X}) and $\operatorname{Conv}(\boldsymbol{Y})$ as

$$D(\boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{\#\text{row}(\boldsymbol{X})} D(\boldsymbol{X}_{i,.}, \boldsymbol{Y}) \text{ and } D(\boldsymbol{X}, \boldsymbol{Y})^{1/2} = \sqrt{D(\boldsymbol{X}, \boldsymbol{Y})}.$$

The number of nonzero entries in a matrix \boldsymbol{A} is denoted as $\|\boldsymbol{A}\|_0$. We define $P_{\ell}(\boldsymbol{H})$ to be the projection of $\boldsymbol{H} \in \mathbb{R}^{k \times n}$ onto the ℓ -sparse set $\{\boldsymbol{X} \in \mathbb{R}^{k \times n} : \|\boldsymbol{X}\|_0 \leq \ell\}$. Moreover, we define the complement as $P_{\ell}^{\perp}(\boldsymbol{H}) = \boldsymbol{H} - P_{\ell}(\boldsymbol{H})$. The support of a matrix $\boldsymbol{H} \in \mathbb{R}^{k \times n}$, $S(\boldsymbol{H})$, is defined as the set of its nonzero coordinates:

$$S({\boldsymbol{H}}) = \{(i,j) \in [k] \times [n] : |{\boldsymbol{H}}_{i,j}| > 0\}.$$

We set $E^{i,j} \in \mathbb{R}^{k \times n}$ to be the matrix with coordinate (i,j) equal to one and other coordinates equal to zero. Throughout this paper, we use $\sigma_{\min}(\boldsymbol{H})$ and $\sigma_{\max}(\boldsymbol{H})$ to denote the smallest and largest singular values of \boldsymbol{H} (respectively). We let $\kappa(\boldsymbol{H}) := \sigma_{\max}(\boldsymbol{H})/\sigma_{\min}(\boldsymbol{H})$ denote the condition number of \boldsymbol{H} . For a convex and subdifferentiable function $f : \mathbb{R}^d \to \mathbb{R}$, $\partial f(\boldsymbol{x})$ denotes the set of subgradients of f at $\boldsymbol{x} \in \mathbb{R}^d$. Proofs of main results have been relegated to the appendix to improve readability.

2 Problem Formulation

Given m data points in \mathbb{R}^n , stacked along the rows of $\mathbf{X} \in \mathbb{R}^{m \times n}$, the goal of AA is to find k archetypes $\mathbf{H}_{1,..}, \dots, \mathbf{H}_{k,..} \in \mathbb{R}^n$ such that: (i) the rows of \mathbf{X} are contained in the convex hull of the rows of \mathbf{H} ; and (ii) the rows of \mathbf{H} are themselves close to the convex hull of the rows of \mathbf{X} . In SAA, we wish to learn a

³Implementation can be found at https://github.com/kayhanbehdin/SparseAA.

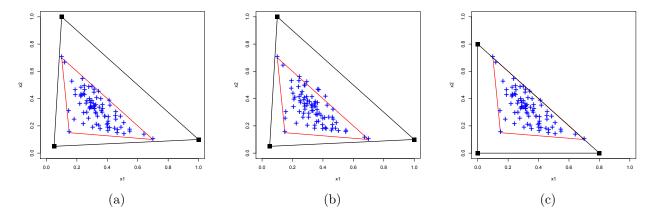


Figure 1: In these figures, blue crosses ('+') represent the data points in \mathbb{R}^2 . We seek to find 3 archetypes such that their convex hull contains the data. Panel (a): the black convex hull (triangle) shows an arbitrary solution to NMF, while the red convex hull shows the exact AA solution that is the smallest triangle containing the data. Panel (b): the black convex hull shows a solution that describes the data with no error, while it is not close to the convex hull of the data. The red convex hull shows the regularized AA solution, which both describes the data well (but with nonzero error) and is close to the data. Panel (c): the black convex hull shows the exact AA solution which does not have any zero coordinate, while the red convex hull shows a solution which is sparse and only has 2 nonzero coordinates. In addition, no other solution with the same sparsity can be found which is closer to the data.

sparse matrix H, i.e., one with few nonzero entries. Equivalently, we seek to learn H such that

$$D(X, H) = D(H, X) = 0 \tag{2}$$

where the first term ensures the data is described by \boldsymbol{H} (as it implies there exists $\boldsymbol{W} \in \mathbb{R}^{m \times k}_{\geq 0}$ such that $\boldsymbol{X} = \boldsymbol{W}\boldsymbol{H}$ and each row of \boldsymbol{W} sums to one). The second term ensures that rows of \boldsymbol{H} are in the convex hull of rows of \boldsymbol{X} . As discussed earlier, the constraint (2) is too restrictive for most practical cases. Javadi and Montanari (2019) propose a relaxed version of (2), where among all the archetypes that contain the data with acceptable accuracy, they choose archetypes that are closest to the convex hull of data. We follow suit, but in addition, we impose a sparsity on \boldsymbol{H} via the constraint $\|\boldsymbol{H}\|_0 \leq \ell$, where, ℓ is a budget on the number of nonzero entries. Specifically, for a pre-specified value of α , we consider:

$$\hat{\boldsymbol{H}} \in \underset{\boldsymbol{H} \in \mathbb{R}_{\geq 0}^{k \times n}}{\operatorname{argmin}} D(\boldsymbol{H}, \boldsymbol{X}) \quad \text{s.t.} \quad D(\boldsymbol{X}_{i,.}, \boldsymbol{H})^{1/2} \leq \alpha, i \in [m]; \quad \|\boldsymbol{H}\|_{0} \leq \ell.$$
(3)

Above the constraint " $D(X_{i,.}, \mathbf{H})^{1/2} \leq \alpha, i \in [m]$ " restricts the data points (rows of \mathbf{X}) to be close to the convex hull of the rows of \mathbf{H} . We also constrain the archetypes (i.e., the right latent factors) to be sparse. Since there are multiple feasible solutions in problem (3), the objective function chooses the archetypes closest to the convex hull of data points, therefore, making the problem more well-defined.

2.1 Model setup

Suppose $X_0 \in \mathbb{R}_{\geq 0}^{m \times n}$ with rank $(X_0) = k$ admits a nonnegative factorization of rank k. That is, there exist $\bar{W} \in \mathbb{R}_{\geq 0}^{m \times k}$, $\bar{H} \in \mathbb{R}_{\geq 0}^{k \times n}$ such that $X_0 = \bar{W}\bar{H}$ and rows of \bar{W} sum to one. This is equivalent to $D(X_0, \bar{H}) = 0$. However, such a factorization is not generally unique. Hence, we let

$$\boldsymbol{H}_0 \in \underset{\boldsymbol{H} \in \mathbb{R}_{\geq 0}^{k \times n}}{\operatorname{argmin}} D(\boldsymbol{H}, \boldsymbol{X}_0) \text{ s.t. } D(\boldsymbol{X}_0, \boldsymbol{H}) = 0$$
 (4)

and assume that \boldsymbol{H}_0 is unique. The choice of \boldsymbol{H}_0 in problem (4) guarantees that \boldsymbol{X}_0 has an exact factorization of the form $\boldsymbol{X}_0 = \boldsymbol{W}_0 \boldsymbol{H}_0$ for some $\boldsymbol{W}_0 \in \mathbb{R}_{\geq 0}^{m \times k}$ such that its rows sum to one and \boldsymbol{H}_0 is defined as in (4). Note that $\operatorname{rank}(\boldsymbol{H}_0) = k$ —otherwise, if $\operatorname{rank}(\boldsymbol{H}_0) < k$, then $\operatorname{rank}(\boldsymbol{X}_0) < k$ which contradicts our assumption on \boldsymbol{X}_0 .

We assume that the observed data matrix X, is given by $X = X_0 + Z$ where Z is additive noise. In what follows, we do not make any distributional assumptions on Z.

Remark 1. Our analysis of robustness is valid without the uniqueness assumption on H_0 . We use uniqueness for simplicity of exposition.

Remark 2. Note that we do not assume that H_0 is sparse, though formulation (3) imposes an explicit cardinality constraint on H. This model misspecification leads to technical challenges: The analysis presented in Javadi and Montanari (2019) does not readily generalize to our setting, and we present a new analysis technique.

2.2 Robustness to noise in Archetypal Analysis

We are interested to see if a solution $\hat{\boldsymbol{H}}$ of (3) is close to \boldsymbol{H}_0 , the underlying set of archetypes. To this end, following Javadi and Montanari (2019), we define a distance between two sets of archetypes $\boldsymbol{H}_1 \in \mathbb{R}^{k_1 \times n}, \boldsymbol{H}_2 \in \mathbb{R}^{k_2 \times n}$ as

$$\mathcal{L}(\boldsymbol{H}_{1}, \boldsymbol{H}_{2}) = \sum_{i=1}^{k_{1}} \min_{j \in [k_{2}]} \|\boldsymbol{H}_{i,.}^{1} - \boldsymbol{H}_{j,.}^{2}\|_{2}^{2}.$$
 (5)

Javadi and Montanari (2019) show that under the so-called *uniqueness* assumption, one has $\mathcal{L}(\boldsymbol{H}_0, \boldsymbol{H}_{\boldsymbol{X}})^{1/2} \leq C \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2$ for some constant C > 0 where $\boldsymbol{H}_{\boldsymbol{X}}$ is the solution to the relaxed AA (i.e. (3) with $\ell = nk$). Note that $\mathcal{L}(\boldsymbol{H}_1, \boldsymbol{H}_2)$ in (5) is a sum of the distances between each row of \boldsymbol{H}_1 and rows of \boldsymbol{H}_2 . Observe that $\mathcal{L}(\boldsymbol{H}_1, \boldsymbol{H}_2)$ is not symmetric in its arguments. In fact, a small value of $\mathcal{L}(\boldsymbol{H}_1, \boldsymbol{H}_2)$ does not imply that $\mathcal{L}(\boldsymbol{H}_2, \boldsymbol{H}_1)$ is also small (see Section 2.3 for details). Definition 1 below presents a formal definition of weak and strong robustness.

Definition 1. (Robustness) An estimator $H \in \mathbb{R}_{>0}^{k \times n}$ is said to be:

- (1) (Weak robustness): Weakly robust if $\mathcal{L}(\boldsymbol{H}_0, \boldsymbol{H})^{1/2} \leq f_1(\max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2)$ where f_1 is an increasing real function that does not depend on \boldsymbol{X} .
- (2) (Strong robustness): Strongly robust if $\mathcal{L}(\boldsymbol{H}, \boldsymbol{H}_0)^{1/2} \leq f_2(\max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2)$ where f_2 is an increasing real function that does not depend on \boldsymbol{X} .

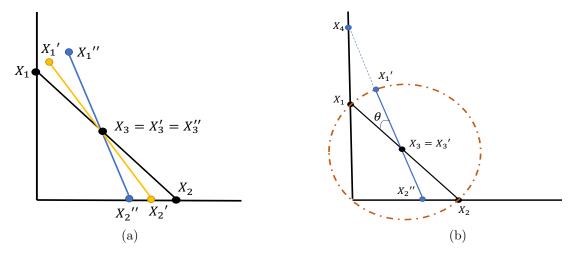


Figure 2: Illustration of Example 1: (a) The noiseless data and two examples of noisy data. (b) Details of the example for a specific value of θ .

Note that based on Definition 1, the result of Javadi and Montanari (2019) is an instance of weak robustness with $f_1(x) = Cx$ for $x \ge 0$ and some C > 0.

2.3 Strong robustness implies weak robustness

Here we explain the differences between weak and strong robustness and provide some intuition around the choice of terminology: strong and weak. Example 1 shows an estimator that is weakly robust but not strongly robust.

Example 1. Let m=3, n=2 and k=2 and

$$m{X}_0 = egin{bmatrix} 0 & 1 \ 1 & 0 \ 1/2 & 1/2 \end{bmatrix}, \quad m{H}_0 = egin{bmatrix} 1 & 0 \ 0 & 1 \end{bmatrix}.$$

For $\theta \in (0, \pi/4)$, we let the noisy data (X_{θ}) and noise matrix (Z_{θ}) be:

$$\boldsymbol{X}_{\theta} = \begin{bmatrix} \sqrt{1 - \cos\theta} \cos(\frac{\pi}{4} - \frac{\theta}{2}) & 1 + \sqrt{1 - \cos\theta} \sin(\frac{\pi}{4} - \frac{\theta}{2}) \\ 1 - \frac{\sin\theta}{\sqrt{2}\sin(\theta + \frac{\pi}{4})} & 0 \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} \quad \text{and}$$

$$\boldsymbol{Z}_{\theta} = \begin{bmatrix} \sqrt{1 - \cos\theta} \cos(\frac{\pi}{4} - \frac{\theta}{2}) & \sqrt{1 - \cos\theta} \sin(\frac{\pi}{4} - \frac{\theta}{2}) \\ -\frac{\sin\theta}{\sqrt{2}\sin(\theta + \frac{\pi}{4})} & 0 \\ 0 & 0 \end{bmatrix}.$$

Figure 2 presents an illustration of Example 1. In this figure (panel (a)), X_1, X_2, X_3 correspond to the original data points (i.e., the rows of X_0) and points X'_1, X'_2, X'_3 and X''_1, X''_2, X''_3 are the noisy data points for two different values of θ . Note that, in all cases, the data points lie on a line. The lines X'_1, X'_2, X'_3 and

 X_1'', X_2'', X_3'' are obtained by rotating the noiseless data line X_1, X_2, X_3 along its center (1/2, 1/2). Figure 2 (panel (b)) shows line X_1, X_2, X_3 and its rotated version X_1', X_2', X_3' (after being rotated by an angle θ). Let

$$m{H}_{ heta} = egin{bmatrix} 0 & \left[1 - rac{\sin heta}{\sqrt{2} \sin(heta + rac{\pi}{4})}
ight] an(heta + rac{\pi}{4}) \ 1 - rac{\sin heta}{\sqrt{2} \sin(heta + rac{\pi}{4})} & 0 \end{bmatrix}.$$

Note that $\boldsymbol{H}_{2,.}^{\theta}$ is the same as the point $\boldsymbol{X}_{2,.}^{\theta}$. The line passing through the noisy data points intersects the y-axis at $\boldsymbol{H}_{1,.}^{\theta}$. For the line, X_{1}', X_{2}', X_{3}' , the point $\boldsymbol{H}_{1,.}^{\theta}$ is given by X_{4} in Figure 2, (b). As a result, $\boldsymbol{X}_{1,.}^{\theta}, \boldsymbol{X}_{2,.}^{\theta}, \boldsymbol{X}_{3,.}^{\theta}$ are on the segment connecting $\boldsymbol{H}_{1,.}^{\theta}$ and $\boldsymbol{H}_{2,.}^{\theta}$ and $\boldsymbol{D}(\boldsymbol{X}_{\theta}, \boldsymbol{H}_{\theta}) = 0$. In addition, $\max_{i \in [3]} \|\boldsymbol{Z}_{i..}^{\theta}\|_{2} \leq \sqrt{2}$ showing the amount of noise added to the data is limited. Moreover,

$$\mathcal{L}(\boldsymbol{H}_0, \boldsymbol{H}_{\theta}) \leq \|\boldsymbol{H}_{1..}^0 - \boldsymbol{H}_{2..}^{\theta}\|_2^2 + \|\boldsymbol{H}_{2..}^0 - \boldsymbol{H}_{2..}^{\theta}\|_2^2 \leq 4$$

for all $\theta \in (0, \pi/4)$, showing \mathbf{H}_{θ} is weakly robust. However, note that

$$\mathcal{L}(\boldsymbol{H}_{\theta}, \boldsymbol{H}_{0}) \geq \min_{i \in [2]} \|\boldsymbol{H}_{1,..}^{\theta} - \boldsymbol{H}_{i,..}^{0}\|_{2}^{2} = \|\boldsymbol{H}_{1,..}^{\theta} - \boldsymbol{H}_{2,..}^{0}\|_{2}^{2}$$
$$= \left(\left[1 - \frac{\sin \theta}{\sqrt{2}\sin(\theta + \pi/4)} \right] \tan(\theta + \pi/4) - 1 \right)^{2}$$

and $\mathcal{L}(\boldsymbol{H}_{\theta}, \boldsymbol{H}_{0}) \to \infty$ as $\theta \to \pi/4$, showing \boldsymbol{H}_{θ} is not strongly robust.

 H_{θ} is not a strongly robust estimator in Example 1 as the first row of H_{θ} can be far from H_0 , the set of underlying archetypes. Weak robustness implies that among recovered archetypes, there is at least one of them which is close to the correct archetypes (for example, the second row of H_{θ} in Example 1 is close to the first row of H_0). On the other hand, strong robustness implies that all recovered archetypes are close to some underlying archetype. This is true because considering the definition of $\mathcal{L}(H_1, H_2)$ in (5), strong robustness limits the distance between each recovered archetype and the set of correct archetypes. However, weak robustness limits the distance between each true archetype and the set of recovered archetypes, therefore, this distance can be small even if some recovered archetypes are far. Theorem 1 below states that strong robustness implies weak robustness (and in light of Example 1, this containment is strict).

Theorem 1 (Strong robustness implies weak robustness). Let us define the quantity

$$b(\mathbf{H}_0) = \max_{i,j \in [k]} \|\mathbf{H}_{i,.}^0 - \mathbf{H}_{j,.}^0\|_2.$$
(6)

Then for any $\boldsymbol{H} \in \mathbb{R}_{\geq 0}^{k \times n}$ we have

$$\mathcal{L}(\boldsymbol{H}_0, \boldsymbol{H}) \le 2kb(\boldsymbol{H}_0)^2 + 2\mathcal{L}(\boldsymbol{H}, \boldsymbol{H}_0). \tag{7}$$

If \boldsymbol{H} is a strongly robust estimator, $\mathcal{L}(\boldsymbol{H}, \boldsymbol{H}_0)$ is bounded and as $b(\boldsymbol{H}_0)$ depends on the underlying archetypes, $b(\boldsymbol{H}_0)$ is finite. Therefore, the right hand side in (7) is bounded and $\mathcal{L}(\boldsymbol{H}_0, \boldsymbol{H})$ is bounded, implying \boldsymbol{H} is also weakly robust.

Theorem 2 shows that a weak/strong robust estimator H serves as a good approximation to the noiseless

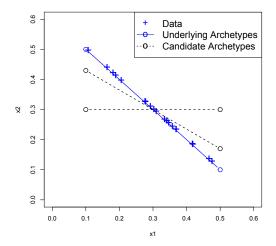


Figure 3: In this figure, blue crosses represent the data points in \mathbb{R}^2 and blue circles and line represent H_0 and its convex hull containing the data. Black circles and lines represent two candidate set of archetypes and their convex hulls. Note that the set that is closer to H_0 describes the data better as anticipated by Theorem 2.

data matrix X_0 .

Theorem 2. For any $H \in \mathbb{R}^{k \times n}_{>0}$ we have

$$D(X_0, H)^{1/2} \le \sqrt{m} \min \left\{ \mathcal{L}(H_0, H)^{1/2}, k \| H_0 \|_F + \mathcal{L}(H, H_0)^{1/2} \right\}.$$
 (8)

If the right hand side in (8) is small, which means \boldsymbol{H} is either a weakly or strongly robust estimator, then the left hand side in (8) is small. This shows that the noiseless data is close to the convex hull of \boldsymbol{H} —so rows of \boldsymbol{H} are good representatives of the noiseless data. See Figure 3 for an illustrative example.

3 Sparse AA

In this section, our primary goal is to show that under certain conditions, both weak and strong robustness hold for the SAA estimator (3).

Theorem 3. Let X_0, H_0 be as described in Section 2.1 and \hat{H} be a solution of problem (3). Set $\alpha = \delta + \beta$ where $\delta = \max_{i \in [m]} \|Z_{i,.}\|_2$ and $\beta = \sqrt{m} \|P_{\ell}^{\perp}(H_0)\|_F$. Moreover, let $X = X_0 + Z$ and $\tilde{X}_0 \in \mathbb{R}^{k \times n}$ be such that

$$ilde{oldsymbol{X}}_{i,.}^0 = \mathop{\mathrm{argmin}}_{oldsymbol{u} \in \{oldsymbol{X}_{j,.}^0: j \in [m]\}} \|oldsymbol{u} - oldsymbol{H}_{i,.}^0\|_2.$$

There exist constants⁴ c_1, \dots, c_{10} depending on $m, n, k, \kappa(\mathbf{H}_0), \sigma_{\min}(\mathbf{H}_0)$ such that the following bounds hold:

1. (Weak Robustness)

$$\mathcal{L}(\boldsymbol{H}_0, \hat{\boldsymbol{H}})^{1/2} \le c_1 D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)^{1/2} + c_2 \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2 + c_3 \|P_{\ell}^{\perp}(\boldsymbol{H}_0)\|_F$$
(9)

⁴To aid readability, the precise expressions for these constants are presented in Section D.3.

2. (Strong Robustness) If

$$c_4 D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)^{1/2} + c_5 \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2 + c_6 \|P_\ell^{\perp}(\boldsymbol{H}_0)\|_F \le c_7,$$
 (10)

we have the following strong robustness guarantee

$$\mathcal{L}(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} \le c_8 D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)^{1/2} + c_9 \max_{i \in [m]} \|\boldsymbol{Z}_{i, \cdot}\|_2 + c_{10} \|P_{\ell}^{\perp}(\boldsymbol{H}_0)\|_F.$$
(11)

In Theorem 3, δ controls the amount of additive noise \mathbf{Z} ; and $\beta = \sqrt{m} \|P_{\ell}^{\perp}(\mathbf{H}_0)\|_F$ captures the sparsity level of \mathbf{H}_0 . If the data is noisy (i.e. δ is large) and/or \mathbf{H}_0 is not sparse (i.e. β is large), the value of α in problem (3) is larger. This implies that $D(\mathbf{X}_{i,.}, \hat{\mathbf{H}})$ for $i \in [m]$ can be potentially larger because of the constraint $D(\mathbf{X}_{i,.}, \mathbf{H})^{1/2} \leq \alpha$ in problem (3)— $\hat{\mathbf{H}}$ might not represent the data points well. However, this is the price we pay to guarantee robustness.

If \mathbf{H}_0 is not ℓ -sparse, or equivalently $\|P_\ell^\perp(\mathbf{H}_0)\|_F > 0$, problem (3) obtains a sparse estimator $\hat{\mathbf{H}}$, which approximates \mathbf{H}_0 . This is an example of model misspecification; and even in this case, Problem (3) leads to an estimator that is weakly and strongly robust. In Theorem 3, the quantity $D(\mathbf{H}_0, \tilde{\mathbf{X}}_0)$ determines how close the underlying model is to the noiseless data; and it depends upon \mathbf{H}_0 and \mathbf{X}_0 . Choosing \mathbf{H}_0 as in (4) results in a smaller value of $D(\mathbf{H}_0, \tilde{\mathbf{X}}_0)$ and improves our bounds. This constant can be zero which is a generalization of the separable case (Donoho and Stodden, 2004), where the underlying archetypes are assumed to be among noiseless data points. In addition, condition (10) ensures that the noise in the data is not too large and the underlying archetypes are suitably sparse — this suffices to derive a strong robustness guarantee for $\hat{\mathbf{H}}$.

In the special case where \mathbf{H}_0 is ℓ -sparse (i.e., $\|P_{\ell}^{\perp}(\mathbf{H}_0)\|_F = 0$) and the underlying model is separable (i.e., $D(\mathbf{H}_0, \tilde{\mathbf{X}}_0) = 0$) the results of Theorem 3 can be simplified as in Corollary 1.

Corollary 1. Let \hat{H} be the solution of problem (3). Under the assumption of Theorem 3 and assuming $||P_{\ell}^{\perp}(\mathbf{H}_0)||_F = D(\mathbf{H}_0, \tilde{\mathbf{X}}_0) = 0$, we have the following:

1. (Weak Robustness)

$$\mathcal{L}(\boldsymbol{H}_{0}, \hat{\boldsymbol{H}})^{1/2} \leq [4mk\kappa(\boldsymbol{H}_{0}) + (1+\sqrt{2})\sqrt{k}(k+k^{3/2})] \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_{2}$$

2. (Strong Robustness) If $[(k + k^{3/2}) + 2mk] \max_{i \in [m]} \|\mathbf{Z}_{i,.}\|_2 \le \sigma_{\min}(\mathbf{H}_0)/(6\sqrt{k}),$

$$\mathcal{L}(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} \leq [7\kappa(\boldsymbol{H}_0)(k+k^{3/2}) + 2(1+\sqrt{2})k^{3/2}m] \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2.$$

Considering that m > k, the results in Corollary 1 can be summarized as

$$\mathcal{L}(\boldsymbol{H}_0, \hat{\boldsymbol{H}})^{1/2} = \mathcal{O}(mk\delta)$$
 and $\mathcal{L}(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} = \mathcal{O}(mk^{3/2}\delta)$

where $\delta = \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2$ (assuming $\kappa(\boldsymbol{H}_0)$ is constant). This shows the bound for the strong robustness quantity $\mathcal{L}(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2}$ is larger. This observation can be explained as follows: strong robustness bounds the distance between each recovered archetype and the underlying archetypes, and as the location of each recovered archetype is unknown and uncertain, the strong robustness quantity deals with more uncertainty

compared to weak robustness and is harder to bound.

Comparison with Javadi and Montanari (2019): Although Javadi and Montanari (2019) do not consider sparsity in their formulation, it is insightful to compare our results with theirs. The results of Javadi and Montanari (2019) are valid under a specific uniqueness assumption on the model—this matches with our assumption in the separable case. Therefore, to compare our results we consider the case $\beta = D(\mathbf{H}_0, \tilde{\mathbf{X}}_0) = 0$ as in Corollary 1. In this regime, the result of Javadi and Montanari (2019) is similar to the first part of the Corollary 1 without any sparsity guarantee on the solution and with different coefficients. They do not provide results similar to the second part of the corollary. Moreover, the bound of Javadi and Montanari (2019) is $\mathcal{L}(\mathbf{H}_0, \hat{\mathbf{H}})^{1/2} = \mathcal{O}(k^{9/4}\delta)$ (assuming other parameters in their bound are constant) which is loose compared to our bound if $m = \mathcal{O}(k^{5/4})$. This shows that our results are tighter when the number of data points is small. Admittedly, the uniqueness assumption of Javadi and Montanari (2019) is more general than the separable case, however, this assumption is difficult to verify except for very simple cases, as discussed by the authors.

3.1 The penalized formulation

Theorem 3 presents robustness guarantees for the constrained SAA problem (3). From an algorithmic viewpoint however, the penalized form:

$$\hat{\boldsymbol{H}}_{\lambda} \in \underset{\boldsymbol{H} \in \mathbb{R}_{>0}^{k \times n}}{\operatorname{argmin}} \quad D(\boldsymbol{X}, \boldsymbol{H}) + \lambda D(\boldsymbol{H}, \boldsymbol{X}) \quad \text{s.t.} \quad \|\boldsymbol{H}\|_{0} \le \ell$$
(12)

is more appealing, and we propose algorithms for this penalized form. In (12), D(X, H) is the data fidelity term, D(H, X) is the regularization term and λ is the regularization parameter. In fact, $\lambda = 0$ is equivalent to setting $\alpha = 0$ in problem (3) (which can lead to an infeasible problem) and $\lambda \to \infty$ is equivalent to removing the data fidelity constraint all together.

We show robustness properties of estimator (12). Proposition 1 establishes both weak (9) and strong (11) robustness. For simplicity, we consider the separable case $(D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0) = 0)$ where \boldsymbol{H}_0 is sparse $(\|\boldsymbol{P}_{\ell}^{\perp}(\boldsymbol{H}_0)\|_F = 0)$.

Proposition 1 (The penalized formulation (12)). Let $\hat{\boldsymbol{H}}_{\lambda}$ be a solution to problem (12). Suppose the assumptions of Theorem 3 hold; and in addition $D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0) = \|P_{\ell}^{\perp}(\boldsymbol{H}_0)\|_F = 0$. There exist constants⁵ $c_{\lambda}^1, c_{\lambda}^2, c_{\lambda}^3$ depending on $m, k, \kappa(\boldsymbol{H}_0), \lambda$ such that $c_{\lambda}^1, c_{\lambda}^2, c_{\lambda}^3 \to \infty$ as $\lambda \to 0$ or $\lambda \to \infty$; and the following holds:

$$\mathcal{L}(\boldsymbol{H}_{0}, \hat{\boldsymbol{H}}_{\lambda})^{1/2} \leq c_{\lambda}^{1} \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_{2}.$$
(13)

Moreover, if

$$c_{\lambda}^{3} \max_{i \in [m]} \| \boldsymbol{Z}_{i,.} \|_{2} \le c_{7},$$
 (14)

then

$$\mathcal{L}(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{H}_{0})^{1/2} \leq c_{\lambda}^{2} \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_{2}.$$
(15)

⁵The values of the constants can be found in the appendix, Section D.4

Note that c_{λ}^{1} , $c_{\lambda}^{2} \to \infty$ as $\lambda \to 0$ or $\lambda \to \infty$ —hence, there is no trivial value of λ that guarantees robustness. In fact, if we have $\lambda = 0$ (this is equivalent to normal NMF), the archetypes need not be close to the data (which is approximated by the underlying archetypes)—they can be far from the underlying archetypes and robustness is not guaranteed. If $\lambda \to \infty$, we do not reduce the recovery error (D(X, H)) which does not result in robustness. This shows the usefulness of using AA as a regularization term.

Remark 3. Results for the penalized form of AA as in Proposition 1 are not discussed in Javadi and Montanari (2019), as far as we can tell.

4 SAA Algorithm

In this section, we propose algorithms to obtain good (feasible) solutions to the penalized SAA problem (12). In Section 4.1, we present a block coordinate descent method and derive its convergence guarantees. As problem (12) is nonconvex, Section 4.2 discusses an initialization scheme based on MIP techniques to keep away from suboptimal solutions. To further improve the quality of the solution obtained from the block coordinate method, we present a heuristic local search algorithm in Section 4.3.

4.1 A Block Coordinate Algorithm

We rewrite problem (12) as follows:

$$\min_{\boldsymbol{W}, \tilde{\boldsymbol{W}}, \boldsymbol{H}} \quad \Psi(\boldsymbol{W}, \tilde{\boldsymbol{W}}, \boldsymbol{H}) := \|\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H}\|_F^2 + \lambda \|\boldsymbol{H} - \tilde{\boldsymbol{W}}\boldsymbol{X}\|_F^2
\text{s.t.} \quad \boldsymbol{H} \ge \mathbf{0}, \boldsymbol{W} \ge 0, \tilde{\boldsymbol{W}} \ge \mathbf{0}
\qquad \boldsymbol{W} \mathbf{1}_k = \mathbf{1}_m, \tilde{\boldsymbol{W}} \mathbf{1}_m = \mathbf{1}_k, \|\boldsymbol{H}\|_0 \le \ell.$$
(16)

We propose a proximal gradient based block coordinate descent algorithm (Xu and Yin, 2017) for (16). We first note that the gradient of the objective function is Lipschitz for every block $(\mathbf{W}, \tilde{\mathbf{W}}, \mathbf{H})$, that is,

$$\begin{split} &\|\nabla_{\boldsymbol{H}}\Psi(\boldsymbol{H}_{1},\boldsymbol{W},\tilde{\boldsymbol{W}}) - \nabla_{\boldsymbol{H}}\Psi(\boldsymbol{H}_{2},\boldsymbol{W},\tilde{\boldsymbol{W}})\|_{F} \leq L_{1}(\boldsymbol{W})\|\boldsymbol{H}_{1} - \boldsymbol{H}_{2}\|_{F}, \\ &\|\nabla_{\boldsymbol{W}}\Psi(\boldsymbol{H},\boldsymbol{W}_{1},\tilde{\boldsymbol{W}}) - \nabla_{\boldsymbol{W}}\Psi(\boldsymbol{H},\boldsymbol{W}_{2},\tilde{\boldsymbol{W}})\|_{F} \leq L_{2}(\boldsymbol{H})\|\boldsymbol{W}_{1} - \boldsymbol{W}_{2}\|_{F} \\ &\|\nabla_{\tilde{\boldsymbol{W}}}\Psi(\boldsymbol{H},\boldsymbol{W},\tilde{\boldsymbol{W}}_{1}) - \nabla_{\tilde{\boldsymbol{W}}}\Psi(\boldsymbol{H},\boldsymbol{W},\tilde{\boldsymbol{W}}_{2})\|_{F} \leq L_{3}(\boldsymbol{X})\|\tilde{\boldsymbol{W}}_{1} - \tilde{\boldsymbol{W}}_{2}\|_{F} \end{split}$$

where $L_1(\mathbf{W}) = 2(\lambda + \sigma_{\max}(\mathbf{W})^2)$, $L_2(\mathbf{H}) = 2\max\{\sigma_{\max}(\mathbf{H})^2, \varepsilon\}$ and $L_3(\mathbf{X}) = 2\lambda\sigma_{\max}(\mathbf{X})^2$ for any fixed $\varepsilon > 0$. Our algorithm follows the block proximal update of Xu and Yin (2017). Specifically, for step size

values $1/2L_1(\boldsymbol{W}_i)$, $1/2L_2(\boldsymbol{H}_i)$, $1/2L_3(\boldsymbol{X})$, at iteration j we perform the following updates:

$$\boldsymbol{H}_{j+1} = \underset{\|\boldsymbol{H}\|_{0} \leq \ell}{\operatorname{argmin}} \left\| \boldsymbol{H} - \left(\boldsymbol{H}_{j} - \frac{1}{2L_{1}(\boldsymbol{W}_{j})} \nabla_{\boldsymbol{H}} \Psi(\boldsymbol{H}_{j}, \boldsymbol{W}_{j}, \tilde{\boldsymbol{W}}_{j}) \right) \right\|_{F}^{2}$$

$$(17)$$

$$\boldsymbol{W}_{j+1} = \underset{\boldsymbol{W} \geq 0}{\operatorname{argmin}} \left\| \boldsymbol{W} - \left(\boldsymbol{W}_{j} - \frac{1}{2L_{2}(\boldsymbol{H}_{j+1})} \nabla_{\boldsymbol{W}} \Psi(\boldsymbol{H}_{j+1}, \boldsymbol{W}_{j}, \tilde{\boldsymbol{W}}_{j}) \right) \right\|_{F}^{2}$$

$$(18)$$

$$\tilde{\boldsymbol{W}}_{j+1} = \underset{\tilde{\boldsymbol{W}} \geq 0}{\operatorname{argmin}} \left\| \tilde{\boldsymbol{W}} - \left(\tilde{\boldsymbol{W}}_{j} - \frac{1}{2L_{3}(\boldsymbol{X})} \nabla_{\tilde{\boldsymbol{W}}} \Psi(\boldsymbol{H}_{j+1}, \boldsymbol{W}_{j+1}, \tilde{\boldsymbol{W}}_{j}) \right) \right\|_{F}^{2}.$$

$$\tilde{\boldsymbol{W}}_{1_{m}=1_{k}}$$

$$(19)$$

After a sweep across the updates (17), (18) and (19) the objective decreases:

$$\Psi(\boldsymbol{H}_{j+1}, \boldsymbol{W}_{j+1}, \tilde{\boldsymbol{W}}_{j+1}) \leq \Psi(\boldsymbol{H}_{j}, \boldsymbol{W}_{j}, \tilde{\boldsymbol{W}}_{j}).$$

Algorithm 1 summarizes the above procedure, where $P_{\text{simplex}}(\boldsymbol{W})$ projects each row of \boldsymbol{W} onto the unit simplex. See Duchi et al. (2008) for an efficient algorithm to calculate P_{simplex} . Before proceeding to the theoretical analysis of Algorithm 1, we need to define stationarity.

Definition 2. We say a point $\theta^* = (H^*, W^*, \tilde{W}^*)$ is stationary for problem (16) if update rules (17), (18) and (19) initialized with θ^* result in the same solution θ^* .

Remark 4. Definition 2 is a generalization of the notion of L-stationarity by Beck and Eldar (2013) to the case of the block proximal method. Moreover, Definition 2 presents a necessary condition for optimality of problem (16).

```
Algorithm 1: SparseAA(H_0, W_0, \tilde{W}_0, \lambda)
```

In Theorem 4, first we show that problem (16) satisfies the convergence conditions of Xu and Yin (2017) and therefore Algorithm 1 converges. Then, we show that the limit point of Algorithm 1 is a stationary point as in Definition 2.

Theorem 4. Suppose $\lambda > 0$ and let $\{(\boldsymbol{H}_j, \boldsymbol{W}_j, \tilde{\boldsymbol{W}}_j)\}_{j \geq 1}$ be the sequence of solutions produced by Algorithm

- 1. The following results hold:
- 1. The sequence $(\boldsymbol{H}_j, \boldsymbol{W}_j, \tilde{\boldsymbol{W}}_j)$ converges to a feasible solution $(\boldsymbol{H}^*, \boldsymbol{W}^*, \tilde{\boldsymbol{W}}^*)$ of (16)
- 2. Let

$$T = \max\{0, \mathbf{H}^* - [1/L_1(\mathbf{W}^*)](-\mathbf{W}^{*T}[\mathbf{X} - \mathbf{W}^*\mathbf{H}^*] + \lambda[\mathbf{H}^* - \tilde{\mathbf{W}}^*\mathbf{X}])\}$$
(20)

and if $\|T\|_0 > \ell$, assume $T_\ell^{\sharp} > T_{\ell+1}^{\sharp}$ where $T_1^{\sharp}, \dots, T_{kn}^{\sharp}$ are entries of T ordered from largest to smallest. Then, the limit point (H^*, W^*, \tilde{W}^*) is a stationary point of (16).

The condition on T above is needed to make sure $P_{\ell}(T)$ is unique. Otherwise, there will be multiple possible solutions for update rule (17) when initialized with θ^* . Note however, that this condition is quite mild, and is unlikely to be violated in practice (due to noise in data).

4.2 Initialization via Mixed Integer Programming (MIP)

Problem (16) is not convex, hence having a good initialization is critical to obtain a high-quality local solution. To initialize Algorithm 1, we use a continuation method as discussed below. We first obtain a solution to (16) for a large value of λ — this leads to the following problem:

$$\min_{\tilde{\boldsymbol{W}} \boldsymbol{H}} \|\boldsymbol{H} - \tilde{\boldsymbol{W}} \boldsymbol{X}\|_F^2 \quad \text{s.t.} \quad \boldsymbol{H} \ge 0; \ \tilde{\boldsymbol{W}} \ge 0; \ \tilde{\boldsymbol{W}} \boldsymbol{1}_m = \boldsymbol{1}_k; \ \|\boldsymbol{H}\|_0 \le \ell. \tag{21}$$

Note that Problem (21) has a convex quadratic objective in \tilde{W} , H and the only source of nonconvexity is the cardinality constraint on H. This is a Mixed Integer Quadratic Problem (MIQP)—while these problems are computationally difficult in the worst-case, recent work (Bertsimas and Van Parys, 2020; Bertsimas et al., 2016; Hazimeh et al., 2020)⁶ has shown that they can be solved to near-optimality using specialized algorithms for large-scale problems. Thusly motivated, we present new algorithms to solve (21) to optimality. Once we obtain a near-optimal solution to (21), we decrease λ and use Algorithm 1 to obtain a feasible solution for (21). We continue this process by successively decreasing λ , and using a solution obtained from the previous (larger) value of λ to initialize Algorithm 1.

To formulate (21) as a MIQP, we first show that the solution of this problem is bounded.

Proposition 2. If (H^*, \tilde{W}^*) is an optimal solution to (21), we have the following bound on H^* :

$$\|\boldsymbol{H}^*\|_F^2 \le k \left(\max_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_2 + \sqrt{k} \min_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_2\right)^2 := b.$$

Based on Proposition 2, we reformulate problem (21) as the following MIQP:

$$\min_{\boldsymbol{H}, \tilde{\boldsymbol{W}}, \boldsymbol{Z}} \quad \|\boldsymbol{H} - \tilde{\boldsymbol{W}} \boldsymbol{X}\|_F^2
\text{s.t.} \quad \boldsymbol{H} \ge 0, \tilde{\boldsymbol{W}} \ge 0, \boldsymbol{Z} \in \{0, 1\}^{k \times n}
\tilde{\boldsymbol{W}} \boldsymbol{1}_m = \boldsymbol{1}_k, \sum_{i,j} \boldsymbol{Z}_{i,j} \le \ell
\boldsymbol{H}_{i,j} \le \sqrt{b} \boldsymbol{Z}_{i,j} \quad \forall (i,j) \in [k] \times [n],$$
(22)

where b is as defined in Proposition 2. Note that the last constraint in (22) does not change the optimal solution because of Proposition 2.

Problem (22) can be formulated and solved (to optimality) by off-the-shelf MIP solvers (e.g., Gurobi, CPLEX, GLPK) for small/moderate instances—however, the runtimes become long as soon as $m, n \sim 100$ or so. With

⁶Note that problem (16) does not admit a MIQP representation, unlike (21).

efficiency in mind, we present a cutting plane approach to obtain a certifiably optimal solution⁷. To this end, we rewrite (22) as a binary convex optimization problem:

$$\min_{\mathbf{Z}} F(\mathbf{Z}) \quad \text{s.t.} \quad \mathbf{Z} \in \{0, 1\}^{k \times n}, \quad \sum_{(i, j) \in [k] \times [n]} \mathbf{Z}_{i, j} \le \ell \tag{23}$$

where, for any $\mathbf{Z} \in [0,1]^{k \times n}$, the objective $F(\mathbf{Z})$ is implicitly defined as the solution to the following convex optimization problem:

$$F(\mathbf{Z}) = \min_{\mathbf{H}, \tilde{\mathbf{W}}} \quad \|\mathbf{H} - \tilde{\mathbf{W}} \mathbf{X}\|_F^2$$
s.t. $\mathbf{H} \ge 0$, $\tilde{\mathbf{W}} \ge 0$, $\tilde{\mathbf{W}} \mathbf{1}_m = \mathbf{1}_k$

$$\mathbf{H}_{i,j} \le \sqrt{b} \mathbf{Z}_{i,j} \quad \forall (i,j) \in [k] \times [n].$$
(24)

Proposition 3 presents some properties of the function $F(\mathbf{Z})$.

Proposition 3. Let $(\boldsymbol{H}^*, \tilde{\boldsymbol{W}}^*)$ be an optimal solution to the minimization problem (24). Then, we have the following:

- 1. The function $\mathbf{Z} \mapsto F(\mathbf{Z})$ is convex on $\mathbf{Z} \in [0,1]^{k \times n}$.
- 2. $G = -\sqrt{b}\Lambda$ is a subgradient of F(Z), where for $(i, j) \in [k] \times [n]$,

$$\boldsymbol{\Lambda}_{i,j} = \begin{cases} 2(\tilde{\boldsymbol{W}}^* \boldsymbol{X} - \boldsymbol{H}^*)_{i,j} & \text{if } (\tilde{\boldsymbol{W}}^* \boldsymbol{X} - \boldsymbol{H}^*)_{i,j} > 0 \\ 0 & \text{if } (\tilde{\boldsymbol{W}}^* \boldsymbol{X} - \boldsymbol{H}^*)_{i,j} \leq 0 \end{cases}.$$

The function $F(\mathbf{Z})$ is convex and subdifferentiable. Specifically, for any $\mathbf{Z}_0 \in \mathbb{R}^{k \times n}$ and any subgradient $\mathbf{G}_0 \in \partial F(\mathbf{Z}_0)$,

$$F(\mathbf{Z}) \ge F(\mathbf{Z}_0) + \langle \mathbf{G}_0, \mathbf{Z} - \mathbf{Z}_0 \rangle. \tag{25}$$

MIP Algorithm: We present an outer approximation algorithm (Duran and Grossmann, 1986) to solve the binary convex program (23). This algorithm starts from an initial point \mathbb{Z}_0 which is feasible for (23). At iteration t, using a list of subgradient-based inequalities (25), we consider the following piecewise linear lower bound of $F(\mathbb{Z})$:

$$F(\mathbf{Z}) \ge \max \left\{ F(\mathbf{Z}_0) + \langle \mathbf{G}_0, \mathbf{Z} - \mathbf{Z}_0 \rangle, \cdots, F(\mathbf{Z}_{t-1}) + \langle \mathbf{G}_{t-1}, \mathbf{Z} - \mathbf{Z}_{t-1} \rangle \right\} \tag{26}$$

where Z_0, \dots, Z_{t-1} are feasible for Problem (23) and $G_t \in \partial F(Z_t)$ for all t. We define Z_t as a minimizer of the right hand side of (26) under the constraints of Problem (23). Mathematically, this can be written as

⁷That is, along with delivering a feasible solution, we also present a dual bound (aka lower bound) on the optimal objective value of (22).

a Mixed Integer Linear Program (MILP)

$$(\boldsymbol{Z}_{t}, \eta_{t}) \in \underset{\boldsymbol{Z}, \eta}{\operatorname{argmin}} \quad \eta$$

$$\text{s.t.} \quad \boldsymbol{Z} \in \{0, 1\}^{k \times n}, \eta \in \mathbb{R}$$

$$\eta \geq F(\boldsymbol{Z}_{i}) + \langle \boldsymbol{G}_{i}, \boldsymbol{Z} - \boldsymbol{Z}_{i} \rangle \quad i = 0, \dots, t - 1$$

$$\sum_{(i, j) \in [k] \times [n]} \boldsymbol{Z}_{i, j} \leq \ell.$$

The optimal objective value of (27) is a lower bound (aka dual bound) for (23); and these lower bounds improve as the iterations progress (i.e., t increases). As the feasible set of problem (23) is finite, after finitely many iterations t, an optimal solution to (23) is found. The optimality gap (OG) of the outer approximation can be calculated as OG = (UB - LB)/UB where LB is the current (and the best) lower bound achieved by the piecewise approximation and UB is the best upper bound for (23) found so far. We summarize the procedure in Algorithm 2 where, 'tol' denotes a pre-specified tolerance level.

Algorithm 2: An outer approximation method to solve (23)

```
t = 1
while OG > tol \ \mathbf{do}
(\mathbf{Z}_t, \eta_t) \text{ are solutions of } (27).
F_{\text{best}} = \min_{i=0,\dots,t-1} F(\mathbf{Z}_i)
OG = (F_{\text{best}} - \eta_t)/F_{\text{best}}
t = t+1
end
```

The optimization Problem (24) is convex in $(\boldsymbol{H}, \tilde{\boldsymbol{W}})$ and we use an accelerated proximal gradient method (Beck and Teboulle, 2009) to solve it.

Note that our proposed algorithm is different from that of Bertsimas and Van Parys (2020) who consider the sparse linear regression problem with an additional ridge regularization. Bertsimas and Van Parys (2020) use an outer approximation algorithm to solve an equivalent convex integer program, with an explicit closed-form expression. In contrast, in our work, the function $F(\mathbf{Z})$ is given (implicitly) by the solution to an optimization problem. Furthermore, our formulation of (22) uses the binary variable as a linear constraint—this is different from Bertsimas and Van Parys (2020) where, the binary variable appears as a nonlinear expression within the objective function.

Numerical results are presented in Section 5.

4.3 Improving Algorithm 1 with Local Search

The block CD method (Algorithm 1) is guaranteed to deliver a stationary point for Problem (16). We present some heuristics to improve the solution quality based on local search, drawing inspiration from the work of Beck and Eldar (2013); Hazimeh and Mazumder (2020) who use local search ideas for a different problem.

Once we are at a stationary point delivered by Algorithm 1, our local search algorithm swaps a coordinate in the support of \boldsymbol{H} with a coordinate from outside the support. That is, a nonzero coordinate of \boldsymbol{H} is set

to zero and a zero coordinate of \boldsymbol{H} is allowed to become nonzero. Then, the optimization is solely done on the coordinate entering the support. If this optimization leads to a lower objective value, we retain the new support. Mathematically, let $(\boldsymbol{H}^*, \boldsymbol{W}^*, \tilde{\boldsymbol{W}}^*)$ be a feasible solution of problem (16) with $\|\boldsymbol{H}^*\|_0 = \ell$. This solution can be an output of Algorithm 1. Suppose $(i_1, j_1) \in S(\boldsymbol{H}^*)$ (here, $S(\boldsymbol{H}^*)$ is the support of \boldsymbol{H}^*) leaves the support and $(i_2, j_2) \notin S(\boldsymbol{H}^*)$ enters the support. We perform an optimization on the coordinate (i_2, j_2) of \boldsymbol{H} to decide whether this change in the support improves the objective function. Let \boldsymbol{E}^{i_1,j_1} be a matrix with all entries equal to zero except coordinate (i_1, j_1) equal to one. We denote $\boldsymbol{H} = \boldsymbol{H}^* - \boldsymbol{H}^*_{i_1,j_1} \boldsymbol{E}^{i_1,j_1}$ as the solution with coordinate (i_1, j_1) removed from the support. The candidate solution with the new support has the form $\boldsymbol{H} + t\boldsymbol{E}^{i_2,j_2}$ for $t \geq 0$. This leads to the following problem:

$$\min_{\boldsymbol{W}, \tilde{\boldsymbol{W}}, t} \|\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H} - t\boldsymbol{W}\boldsymbol{E}^{i_2, j_2}\|_F^2 + \lambda \|\boldsymbol{H} + t\boldsymbol{E}^{i_2, j_2} - \tilde{\boldsymbol{W}}\boldsymbol{X}\|_F^2$$
s.t. $t \ge 0, \boldsymbol{W} \ge 0, \tilde{\boldsymbol{W}} \ge 0$

$$\boldsymbol{W} \mathbf{1}_k = \mathbf{1}_m, \tilde{\boldsymbol{W}} \mathbf{1}_m = \mathbf{1}_k$$
(28)

where, we are optimizing over $(\mathbf{W}, \tilde{\mathbf{W}}, t)$ for a given (i_2, j_2) and (i_1, j_1) . For a fixed value of t, the optimal values of \mathbf{W} and $\tilde{\mathbf{W}}$ in (28) are given as

$$\hat{\boldsymbol{W}} \in \underset{\boldsymbol{W}}{\operatorname{argmin}} \|\boldsymbol{X} - \boldsymbol{W}(\boldsymbol{H} + t\boldsymbol{E}^{i_2, j_2})\|_F^2 \quad \text{s.t. } \boldsymbol{W} \ge 0; \quad \boldsymbol{W} \boldsymbol{1}_k = \boldsymbol{1}_m$$
 (29)

$$\hat{\tilde{\boldsymbol{W}}} \in \underset{\tilde{\boldsymbol{W}}}{\operatorname{argmin}} \quad \|\boldsymbol{H} + t\boldsymbol{E}^{i_2, j_2} - \tilde{\boldsymbol{W}}\boldsymbol{X}\|_F^2 \quad \text{ s.t. } \quad \tilde{\boldsymbol{W}} \ge 0; \quad \tilde{\boldsymbol{W}}\boldsymbol{1}_m = \boldsymbol{1}_k.$$
(30)

Problems (29) and (30) are convex and can be efficiently solved by standard first order methods such as proximal gradient. Note that these first order methods also benefit from warm-starts available from prior estimates of $(\mathbf{W}, \tilde{\mathbf{W}})$. Once \mathbf{W} and $\tilde{\mathbf{W}}$ are updated by (29) and (30), the value of t that minimizes (28) with $\mathbf{W} = \hat{\mathbf{W}}$ and $\tilde{\mathbf{W}} = \hat{\tilde{\mathbf{W}}}$ is:

$$t = \max \left\{ \frac{\sum_{r=1}^{m} U_{r,j_2} W_{r,i_2} - \lambda V_{i_2,j_2}}{\lambda + \|W_{.,i_2}\|_2^2}, 0 \right\}$$
(31)

where U = X - WH and $V = H - \tilde{W}X$.

We use an alternating optimization scheme where the three updates (29),(30) and (31) are performed sequentially until convergence. These updates result in a descent method by construction, though there may not be a strict decrease in the objective value (in which case, the swap may not result in a better solution). In the discussion above, we assumed a fixed pair of indices (i_1, j_1) and (i_2, j_2) . Ideally, we would like to try all possible choices of such indices and consider the one that leads to the maximal decrease in objective value (if any). As this is computationally intensive, we use a heuristic to select a suitable pair of indices. We choose (i_1, j_1) to be the smallest nonzero entry in \mathbf{H}^* :

$$(i_1, j_1) \in \underset{(i,j) \in S(\boldsymbol{H}^*)}{\operatorname{argmin}} \boldsymbol{H}_{i,j}^*. \tag{32}$$

For the pair (i_2, j_2) from outside the current support of H, we choose the coordinate of H^* that has the

Algorithm 3: A Local Search improvement for Algorithm 1

```
Initialize with \boldsymbol{H}^*, \boldsymbol{W}^*, \tilde{\boldsymbol{W}}^*.

while not converged do

if \|\boldsymbol{H}^*\|_0 < \ell then

Choose (i_2, j_2) as in (33).

\boldsymbol{H} = \boldsymbol{H}^*
end
else

Choose (i_1, j_1) as in (32) and (i_2, j_2) as in (33).

\boldsymbol{H} = \boldsymbol{H}^* - \boldsymbol{H}^*_{i_1, j_1} \boldsymbol{E}^{i_1, j_1}
end
while not converged do

Update \boldsymbol{W}, \tilde{\boldsymbol{W}} and t via (29), (30) and (31).
end

if \Psi(\boldsymbol{H} + t\boldsymbol{E}^{i_2, j_2}, \boldsymbol{W}, \tilde{\boldsymbol{W}}) < \Psi(\boldsymbol{H}^*, \boldsymbol{W}^*, \tilde{\boldsymbol{W}}^*) then

\boldsymbol{H}^* = \boldsymbol{H} + t\boldsymbol{E}^{i_2, j_2}, \boldsymbol{W}^* = \boldsymbol{W}, \tilde{\boldsymbol{W}}^* = \tilde{\boldsymbol{W}}
end
end
```

smallest (most negative) gradient of the objective function:

$$(i_2, j_2) \in \underset{(i,j) \in [k] \times [n] \setminus S(\boldsymbol{H}^*)}{\operatorname{argmin}} \frac{\partial}{\partial \boldsymbol{H}_{i,j}} \Psi(\boldsymbol{H}^*, \boldsymbol{W}^*, \tilde{\boldsymbol{W}}^*).$$
 (33)

The overall procedure of local search is shown in Algorithm 3.

5 Numerical Experiments

In this section, we discuss results of our numerical experiments on synthetic and real data and investigate how our framework performs in practice. Our experiments are done on a computer equipped with Intel(R) Core(R) i7 6700HQ CPU @ 2.60GHz, running Microsoft(R) Windows(R) 10 and using 16GB of RAM. We have implemented all algorithms in Julia and we use Gurobi(R) to solve MILPs arising in our initialization scheme. An implementation of our framework in Julia is available at:

https://github.com/kayhanbehdin/SparseAA.

5.1 Synthetic Data

In this section, we consider synthetic data to validate the theory we developed in Sections 3 and 4; and gather further insights into the operating characteristics of SAA.

Dataset generation: The entries of \mathbf{H}_0 are drawn iid from Unif[0,1] and 20% of entries are set to zero at random to produce a matrix with at most 0.8nk nonzero entries. Independent of \mathbf{H}_0 , the entries of \mathbf{W}_0 are drawn iid from Unif[0,1] and each row is normalized to sum to one. The noiseless data is $\mathbf{X}_0 = \mathbf{W}_0 \mathbf{H}_0$ and the noisy data is produced as $\mathbf{X} = \max\{\mathbf{X}_0 + \mathbf{Z}, 0\}$ where entries of \mathbf{Z} are from an independent Gaussian ensemble with mean zero and variance σ_z^2 .

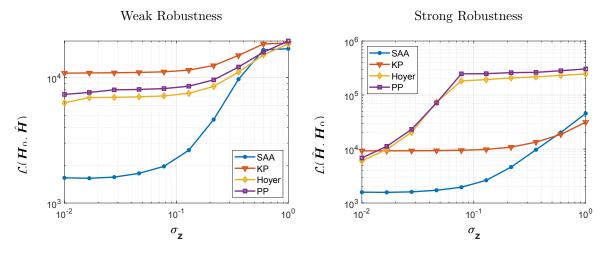


Figure 4: Figure comparing our method (SAA) with three other sparse NMF methods. We compare weak (left panel) and strong (right panel) robustness for varying values of σ_z on the synthetic data in Section 5.1.

5.1.1 Understanding robustness

We compare the performance of our algorithm with other sparse NMF algorithms in terms of robustness of the solution. We consider algorithms by Kim and Park (2007) (shown as KP), Peharz and Pernkopf (2012) (shown as PP) and Hoyer (2004) (shown as Hoyer). We set $m = 200, k = 15, n = 5000, \lambda = 1$. We consider two settings: First, we keep the sparsity level ℓ fixed and change the noise level σ_z ; Second, we keep σ_z fixed and vary the sparsity level ℓ .

Robustness versus varying σ_z : First, we set $\ell/nk = 0.5$ and tune parameters for different algorithms to get solutions with 0.5nk nonzeros. Specifically, as KP considers an ℓ_1 regularized version of NMF, we start with a small value of the ℓ_1 regularization parameter and gradually increase it till we reach the target sparsity-level. PP uses an ℓ_0 constrained version of original NMF (without archetypal regularization) so we set the ℓ_0 sparsity level to ℓ . The sparsity constraint of Hoyer is set such that the result has ℓ nonzeros. We vary σ_z and plot the average value of weak $(\mathcal{L}(H_0, H))$ and strong $(\mathcal{L}(H, H_0))$ robustness quantities. The results for this scenario are shown in Figure 4. As it can be seen, SAA almost always outperforms other algorithms in terms of strong and weak robustness of solutions. Specifically, the difference between SAA and other algorithms is most noticeable when the noise is small. This is expected as other algorithms in our experiments do not use any regularization that results in robustness. Moreover, solutions of SAA become less robust as noise is increased, as anticipated by Theorem 3 and Proposition 1. In addition, it is interesting to note that the weak robustness quantity in Figure 4 [left panel] is smaller than the strong robustness quantity in Figure 4 [right panel]. This is expected based on our discussions in Section 3.

Robustness versus varying ℓ : To compare the performance of different algorithms for varying values of ℓ , we do another set of experiments. We consider a setup similar to the previous experiment. However, we fix $\sigma_z = 0.1$ and change the value of ℓ (while keeping the underlying model sparsity fixed at 0.8nk). This shows how well different algorithms can deal with misspecification in the underlying model. The results for this case are shown in Figure 5. We see that SAA outperforms other algorithms and provides the most robust solutions (both in terms of weak and strong robustness). In addition, we observe that as ℓ is decreased, the

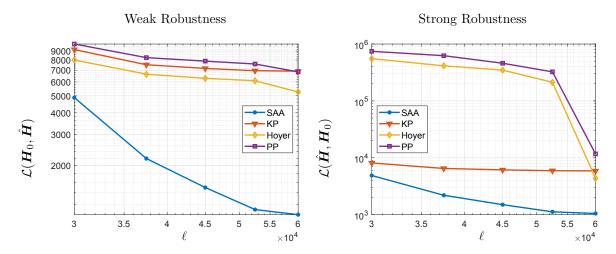


Figure 5: Comparison of robustness quantities for varying values of ℓ on the synthetic data in Section 5.1.

solution becomes less robust as anticipated by Theorem 3.

5.1.2 Usefulness of MIP-based Initialization and Local Search

We perform numerical experiments to show the usefulness of the initialization procedure in Section 4.2 and the local search scheme in Section 4.3. We show these algorithms improve upon a baseline initialization. We set m = 200, k = 20, n = 12000 and $\lambda = 1$ and vary σ_z and ℓ (see Table 1). We consider three cases:

- (i) Algorithm 1 initialized with $\mathbf{H} = \mathbf{0}$ (shown as Zero in Table 1).
- (ii) Algorithm 1 with MIP initialization (see Section 4.2) and using warm-start continuation over 8 values of λ on a logarithmic scale from 30 to 1. (Shown as SAA in Table 1).
- (iii) Improving solution from (ii) with local search discussed in Section 4.3. This is denoted by SAA+LS in Table 1.

Algorithm 2 is limited to 20 minutes of maximum runtime and the best solution is returned. The average final cost function achieved by three methods explained above (over 5 independently generated datasets) are reported in Table 1. As it can be seen, our initialization scheme achieves a significantly lower objective value compared to the baseline (Zero) while being computationally feasible for such data size. In addition, our local search algorithm can improve the objective value as well as the support in a reasonable time (the number of changes in support for SAA+LS is reported in the parentheses). In our experiments, Algorithm 1 terminates in less than 15 seconds.

5.2 The Face Dataset (Samaria and Harter, 1994)

A classical application of sparse NMF is in face detection and recognition (Hoyer, 2004). The goal is to obtain a low-rank representation of a dataset of human faces under different lighting and shadow conditions and also different angles of photography. Hoyer (2004) show the effect of sparsity in finding such representations of the data. In particular, they show that sparse NMF leads to part-based representations where each factor represents one part of the face. Here, we are interested in finding the effect of AA as well as the combined effect of AA and sparsity. We use the AT&T database of faces (Samaria and Harter, 1994) which consists

σ_z	Method	$\ell/nk = 0.5$	$\ell/nk = 0.65$	$\ell/nk = 0.8$
	Zero	47785	42123	37566
0.01	SAA	32193	29766	27139
	SAA+LS	31867 (12.8)	29729 (13.4)	27091 (20.2)
0.1	Zero	69366	63776	59276
	SAA	54784	52253	49228
	SAA+LS	54465 (18.4)	52159 (19.8)	49131 (24)
0.5	Zero	448088	440974	434590
	SAA	426238	420134	416901
	SAA+LS	424776 (22.6)	419932 (26.8)	416628 (32)

Table 1: Comparison of zero initialization, SAA and SAA+LS in Section 5.1. The number in parentheses for SAA+LS shows the number of changes in the support after local search.

of 40 different people and 10 different photos of each person, 400 images in total. Each image is a grayscale 92×112 image, which is converted to a vector of length 10304. We then concatenate these 400 vectorized images into one matrix of size 400×10304 matrix. We consider k=25 (following Hoyer (2004)) for this dataset and do the factorization based on problem (12) for different values of λ and ℓ . The estimated representations of the data (rows of resulting \boldsymbol{H} which are reshaped into 92×112 images) are shown in Figure 6. We use the MIP initialization and continuation framework (over 8 values of λ) in Section 4.2. In the rest of this section, we discuss the connections between the robustness theory we developed and our numerical results.

We first explore the difference between AA and basic NMF. By comparing Figures 6 (a), (b), we deduce that as λ is increased, the resulting factors appear to become more similar to human faces—making it easier to recognize the people in the dataset. As λ is decreased, the results become more abstract and the images do not resemble human faces anymore. Considering that the theory we developed in Section 3 holds under mild assumptions, the representation achieved from Problem (12) with $\lambda = 0.4$ more likely corresponds to a robust solution—that is, the solution is closer to the underlying representation (we do not expect $\lambda = 0$ to result in robustness). Intuitively, we expect the underlying model that produces the face images to resemble the people in the dataset—this suggests that the solution achieved by $\lambda = 0.4$ is more robust.

As discussed by Hoyer (2004), adding sparsity to NMF produces part-based representations of the data. This can be also seen in our experimental results in Figure 6. By forcing the solution H to be sparse, we notice that each set of solutions (in Figures 6 (c), (d)) includes two groups of factors. The first group consists of complete faces (columns and rows 2 to 4 in each set) and the second group consists of parts of a face (the border factors). In fact, this can be interpreted as each face being a combination of an overall shape of a human face and additional details arising from different parts of the face. The factors that contain a complete face in Figures 6 (c), (d) represent the overall shape of a human face, while other factors represent different parts of a face, like the forehead, cheeks, eyes and also the background of images. Once again, considering that our theoretical development in Section 3 is valid in the sparse NMF case, we expect the factors recovered by $\lambda = 0.4$ to be more robust. Our hypothesis is the solution with $\lambda = 0.4$ is more robust as factors in Figure 6 (d) appear to more closely resemble human faces (compare the central columns in Figures 6 (c), (d)).

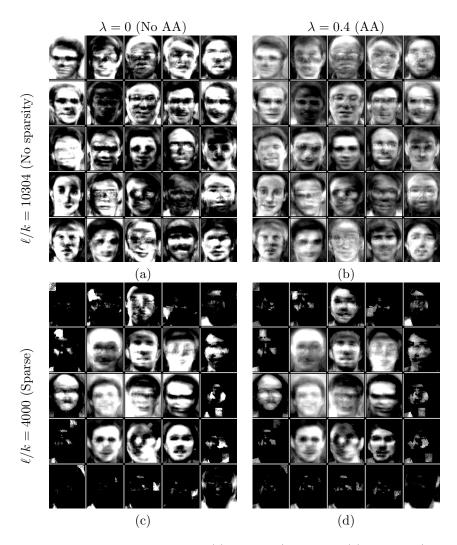


Figure 6: The resulting face images in Section 5.2: (a) $\lambda = 0$, $\ell/k = 10304$ (b) $\lambda = .4$, $\ell/k = 10304$ (c) $\lambda = 0$, $\ell/k = 4000$ (d) $\lambda = .4$, $\ell/k = 4000$

5.3 Cancer Gene Expression Example (Ramaswamy et al., 2001)

It is well-known that a primary goal of AA is to find a few representative points for a collection of data points. These representative points are useful in cluster analysis, where data points are put into a few clusters based on a suitable similarity/dissimilarity measure. In AA, each archetype (i.e., a row of \mathbf{H}) can be considered as a cluster center and data points are assigned to clusters based on their proximity to different archetypes. As a result, each row of matrix \mathbf{H} (or each archetype) is considered as a center and each data point is assigned to the closest row of \mathbf{H} (Mørup and Hansen, 2012).

An important application area of sparse NMF for clustering is in computational biology. Specifically, this problem has been considered by Kim and Park (2007) where the authors provide biological interpretations of sparsity in the context of NMF and do an extensive analysis of sparse NMF for microarray data. Here, we are interested in the clustering performance of our method. To this end, we consider a real dataset: the 14 Cancers Gene Expression dataset (Ramaswamy et al., 2001). This dataset consists of gene expression

data of 198 samples and 14 different types of tumors. There are 16,063 features in the dataset, however, the data is not nonnegative. Therefore, we use a trick introduced by Kim and Tidor (2003) to transform the data: each feature is divided into two new features where one contains nonnegative coordinates (and zero elsewhere) and the other one contains the absolute value of negative coordinates (and zero elsewhere). Consequently, this leads to 32,126 nonnegative features—the data matrix is given by $X \in \mathbb{R}^{198 \times 32126}_{\geq 0}$. The rank of the factorization is 14, the number of different types of tumors in the dataset. The *i*-th data point $(i \in [198])$ belongs to the cluster j_i :

$$j_i = \underset{j \in [14]}{\operatorname{argmin}} \| \boldsymbol{X}_{i,.} - \boldsymbol{H}_{j,.} \|_2$$

where H is the resulting matrix of archetypes. To compare the performance of different algorithms (discussed below), we use two metrics, Purity and Entropy (Kim and Park, 2008). Let for $i \in [198]$, $j_i^* \in [14]$ denote the true cluster of point i and $j_i \in [14]$ denote the estimated cluster for the same point. Let m_r^u be the number of samples that belong to the true cluster u but are estimated to be in cluster r. Equivalently,

$$m_r^u = |\{i \in [m] : j_i^* = u, j_i = r\}|.$$

The metrics Purity and Entropy are defined as

$$\text{Purity} = \frac{1}{m} \sum_{r=1}^k \max_{u \in [k]} m_r^u \quad \text{and} \quad \text{Entropy} = -\frac{1}{m \log_2 k} \sum_{r=1}^k \sum_{u=1}^k m_r^u \log_2 \frac{m_r^u}{m_r}$$

where $m_r = \sum_{u=1}^k m_r^u$. A larger value of Purity and a smaller value of Entropy imply a better clustering performance.

The results of clustering performance of different algorithms is reported in Table 2. SAA is our proposed framework (we use continuation over 8 values of λ from 30 to 1 as in Section 4.2), AA is the algorithm proposed by Javadi and Montanari (2019) which does not enforce any sparsity. KP and PP are as introduced before. We also include Kmeans in our experiments as a baseline for the clustering performance. Hoyer did not provide interpretable results on this dataset and therefore is not included in the table. As it can be seen, among algorithms that enforce sparsity, SAA performs the best in terms of clustering. In fact, SAA is at par with Kmeans, while providing a solution that is two times more sparse. AA has the best clustering performance, however, it fails to provide a sparse solution. Other algorithms provide sparse solutions, but their clustering performance is not as good as SAA. In our experiments in this section, all NMF-based methods terminated in less than a minute.

	SAA	AA	Kmeans	KP	PP
$\ \boldsymbol{H}\ _0/kn$	0.350	0.649	0.710	0.365	0.385
Purity	0.660	0.868	0.654	0.446	0.477
Entropy	0.361	0.216	0.375	0.723	0.690

Table 2: Performance of different algorithms for the gene expression dataset in Section 5.3

5.4 Scene Categorization dataset (Xiao et al., 2010)

Finally, we consider another popular application of AA arising in the context of image categorization (Abrol and Sharma, 2020; Chen et al., 2014). Given a collection of photos, AA can be used to identify a small subset of photos as their representatives (archetypes). Based on the mathematical formulation of AA, archetypes are expected to represent extreme scenes and objects present in the data, for instance, they should categorize different indoor/outdoor settings or city/nature scenes.

As far as we can tell, NMF with an explicit ℓ_0 regularization has not been used before to address the problem of scene categorization. However, in view of Theorem 3, we do not anticipate that additional sparsity will reduce the robustness properties of the estimator if the underlying matrix of archetypes is sparse. In addition, a sparse model may be desirable in terms of compressed storage. We apply SAA on the Scene Categorization dataset (Xiao et al., 2010). We select 12 different scenes that consist of different indoor and outdoor settings (2617 images in total). These scenes are toll plaza, hospital exterior, harbor, electricity station, underwater, youth hostel, valley, ski resort, football stadium, residential neighborhood, vineyard and iceberg. We extract and concatenate GIST and HOG features (Xiao et al., 2010) and implement different sparse NMF algorithms on the data with k = 12. As estimated archetypes in the feature space cannot be visualized, we use the closest data point to each archetype to visualize the result.

First, we consider AA and SAA with $\ell/kn=0.5$ (we use the continuation framework in Section 4.2 and choose the value of λ that maximizes purity). The resulting visualization of archetypes for these two cases is the same. The visualization of estimated archetypes is shown in Figure 7 (a). We observe that the resulting archetypes appear to span the 12 different scenes in the dataset. Figure 7 (b) shows the resulting visualization for PP where the resulting archetypes matrix is set to have at most 0.65nk nonzeros. As it can be seen, PP can identify 10 distinct scenes and chooses the electricity station and the toll plaza twice. Figure 7 (c) shows the results for KP with the same sparsity as PP. This algorithm only identifies 5 distinct scenes. In summary, our SAA algorithm works as well as AA in terms of identifying different scenes while providing a sparse solution. This shows a sparse solution can be achieved without losing categorization performance.

6 Conclusion

In this paper, we consider the problem of sparse NMF with archetypal regularization where the goal is to represent a collection of data points as nonnegative linear combinations of a few nonnegative sparse factors. Javadi and Montanari (2019) recently showed that NMF (without sparsity) with archetypal regularization leads to robustness—factors learnt from noisy data are close to the underlying factors that generate the noiseless data. We generalize the notion of robustness to (a) strong robustness that implies each estimated archetype is close to the underlying archetypes and (b) weak robustness that implies there exists at least one recovered archetype that is close to the underlying archetypes. Javadi and Montanari (2019) is an instance of the notion of weak robustness presented herein. We show that under minimal assumptions, robustness in sparse NMF can be achieved by considering a sparsity constrained regularized AA problem, even if the underlying archetypes are not sparse. We present a block coordinate algorithm to get a good solution to the sparse AA problem and also an initialization framework using mixed integer programming that leads to better numerical results. We also present a local search algorithm that improves the quality of the solution of

SAA (this paper)



PP (Peharz and Pernkopf, 2012)



KP (Kim and Park, 2007)



Figure 7: The visualization of archetypes achieved by different algorithm for the scene categorization data in Section 5.4: (a) SAA (b) PP (c) KP.

our block coordinate algorithm. Numerical experiments on synthetic and real datasets shed further insights into the theoretical developments pursued in this paper.

7 Acknowledgements

Rahul Mazumder would like to thank Kushal Dey for helpful discussions in the initial stages of this work. This research was partially supported by grants from the Office of Naval Research (ONR-N000141812298), National Science Foundation (NSF-IIS-1718258), IBM and Liberty Mutual Insurance.

Proofs and Technical Details

A Additional Notation

We define

$$\tilde{D}(\boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{\#\text{row}(\boldsymbol{X})} \sqrt{D(\boldsymbol{X}_{i,.}, \boldsymbol{Y})} \text{ and } \tilde{\mathcal{L}}(\boldsymbol{X}, \boldsymbol{Y}) = \sum_{i=1}^{\#\text{row}(\boldsymbol{X})} \min_{j \in [\#\text{row}(\boldsymbol{Y})]} \|\boldsymbol{X}_{i,.} - \boldsymbol{Y}_{j,.}\|_{2}.$$
(34)

For a set S, S^c denotes the complement of the set. We use \mathbb{I}_k to denote the identity matrix of size k.

B Technical Details of The Toy Example

For Figure 1 Panel (a), we let

$$m{H}_0 = egin{bmatrix} 0.15 & 0.15 \\ 0.1 & 0.7 \\ 0.7 & 0.1 \end{bmatrix}$$

and produce 50 data points. To do so, each entry of W_0 is drawn from an independent uniform distribution in [0,1] and each row is normalized to sum to one. The noiseless data matrix is $X_0 = W_0 H_0$ and we add three rows of H_0 to this to a obtain separable problem. In this case, we can see $D(X_0, H_0) = D(H_0, X_0) = 0$ which is the exact AA solution of Cutler and Breiman (1994). The red convex hull is $Conv(H_0)$. Let

$$m{H}_1 = egin{bmatrix} 0.05 & 0.05 \\ 1 & 0.1 \\ 0.1 & 1 \end{bmatrix}.$$

The black convex hull is $Conv(\mathbf{H}_1)$ for which we have $D(\mathbf{X}_0, \mathbf{H}_1) = 0$ but $D(\mathbf{H}_1, \mathbf{X}_0) > 0$.

For Figure 1 Panel (b), the data X is produced by $X = X_0 + Z$ where Z has zero-mean iid normal coordinates with variance of 0.1. In this case, we have for every data point i, $D(X_{i,.}, H_0) \leq 0.1$ making H_0 a feasible solution for the regularized AA of Javadi and Montanari (2019) and as $D(H_0, X) = 0$, this is the solution of the regularized AA problem (the red convex hull). However, we have $D(X, H_1) = 0$ and $D(H_1, X) > 0$, so the black convex hull is not an optimal solution. Finally, let

$$m{H}_2 = egin{bmatrix} 0 & 0 \ 0 & 0.8 \ 0.8 & 0 \end{bmatrix}.$$

In Figure 1 Panel (c), the red convex hull is $Conv(\mathbf{H}_2)$ and the black one is $Conv(\mathbf{H}_0)$. We have $\|\mathbf{H}_2\|_0 = 2$ and $\|\mathbf{H}_0\|_0 = 6$ so the black convex hull is not sparse. In addition, among all solutions that have $\|\mathbf{H}\|_0 = 2$ and $D(\mathbf{X}_0, \mathbf{H}) = 0$, the quantity $D(\mathbf{H}, \mathbf{X}_0)$ is minimized for the red convex hull, making it the sparse archetypal solution.

C Technical Lemmas

Lemma 1. For any two matrices $X \in \mathbb{R}^{m_1 \times n}$ and $Y \in \mathbb{R}^{m_2 \times n}$, we have:

$$\frac{1}{\sqrt{m_1}}\tilde{D}(\boldsymbol{X}, \boldsymbol{Y}) \le D(\boldsymbol{X}, \boldsymbol{Y})^{1/2} \le \tilde{D}(\boldsymbol{X}, \boldsymbol{Y})$$
(35)

$$\frac{1}{\sqrt{m_1}}\tilde{\mathcal{L}}(\boldsymbol{X},\boldsymbol{Y}) \le \mathcal{L}(\boldsymbol{X},\boldsymbol{Y})^{1/2} \le \tilde{\mathcal{L}}(\boldsymbol{X},\boldsymbol{Y}), \tag{36}$$

where, recall \tilde{D} and $\tilde{\mathcal{L}}$ are as defined in (34).

Proof. Define the vector $\mathbf{u} \in \mathbb{R}^{m_1}$ such that $\mathbf{u}_i = \sqrt{D(\mathbf{X}_{i,.}, \mathbf{Y})}$. Note that

$$\frac{1}{\sqrt{m_1}}\tilde{D}(\boldsymbol{X},\boldsymbol{Y}) = \frac{1}{\sqrt{m_1}}\|\boldsymbol{u}\|_1 \leq \|\boldsymbol{u}\|_2 = D(\boldsymbol{X},\boldsymbol{Y})^{1/2} \leq \|\boldsymbol{u}\|_1 = \tilde{D}(\boldsymbol{X},\boldsymbol{Y})$$

which establishes (35). The proof of (36) is similar.

Lemma 2. If $D(A_1, A_2) \leq D(B_1, B_2)$ where $A_1, A_2 \in \mathbb{R}^{m \times n}, B_1, B_2 \in \mathbb{R}^{k \times n}$, we have

$$\tilde{D}(\boldsymbol{A}_1, \boldsymbol{A}_2) \leq \sqrt{m} \tilde{D}(\boldsymbol{B}_1, \boldsymbol{B}_2).$$

Proof. We make use of Lemma 1. The proof follows from:

$$\tilde{D}(A_1, A_2) \le \sqrt{m}D(A_1, A_2)^{1/2} \le \sqrt{m}D(B_1, B_2)^{1/2} \le \sqrt{m}\tilde{D}(B_1, B_2).$$

Lemma 3. Suppose $X \in \mathbb{R}^{m_1 \times n}, Y \in \mathbb{R}^{m_2 \times n}, Z \in \mathbb{R}^{m_3 \times n}$, then

$$\tilde{D}(\boldsymbol{X}, \boldsymbol{Z}) \leq \tilde{\mathcal{L}}(\boldsymbol{X}, \boldsymbol{Y}) + m_1 \tilde{D}(\boldsymbol{Y}, \boldsymbol{Z}).$$

Proof. Fix $i \in [m_1]$. For any $\mathbf{u} \in \mathbb{R}^n$, we have

$$\sqrt{D(\boldsymbol{X}_{i,.}, \boldsymbol{Z})} = \min_{\boldsymbol{v} \in \text{Conv}(\boldsymbol{Z})} \|\boldsymbol{X}_{i,.} - \boldsymbol{v}\|_{2}$$

$$\leq \min_{\boldsymbol{v} \in \text{Conv}(\boldsymbol{Z})} \{ \|\boldsymbol{X}_{i,.} - \boldsymbol{u}\|_{2} + \|\boldsymbol{u} - \boldsymbol{v}\|_{2} \}$$

$$= \|\boldsymbol{X}_{i,.} - \boldsymbol{u}\|_{2} + \min_{\boldsymbol{v} \in \text{Conv}(\boldsymbol{Z})} \|\boldsymbol{u} - \boldsymbol{v}\|_{2}.$$
(37)

Let

$$\boldsymbol{u} = \operatorname*{argmin}_{\boldsymbol{p} \in \{\boldsymbol{Y}_{j,.}: j \in [\#\text{row}(\boldsymbol{Y})]\}} \|\boldsymbol{X}_{i,.} - \boldsymbol{p}\|_2.$$

As a result, noting that $\tilde{D}(Y, Z) = \sum_{j=1}^{m_2} \sqrt{D(Y_{j,.}, Z)}$ and u is a row of Y,

$$\min_{\boldsymbol{v} \in \text{Conv}(\boldsymbol{Z})} \|\boldsymbol{u} - \boldsymbol{v}\|_{2} = \sqrt{D(\boldsymbol{u}, \boldsymbol{Z})} \le \sum_{j=1}^{m_{2}} \sqrt{D(\boldsymbol{Y}_{j,.}, \boldsymbol{Z})} = \tilde{D}(\boldsymbol{Y}, \boldsymbol{Z}).$$
(38)

Using the definition of u and (38) we can bound the rhs in (37) to get:

$$\sqrt{D(\boldsymbol{X}_{i,.},\boldsymbol{Z})} \leq \min_{j \in m_2} \|\boldsymbol{X}_{i,.} - \boldsymbol{Y}_{j,.}\|_2 + \tilde{D}(\boldsymbol{Y},\boldsymbol{Z}).$$

By summing the above over i, we have:

$$\begin{split} \tilde{D}(\boldsymbol{X}, \boldsymbol{Z}) &= \sum_{i=1}^{m_1} \sqrt{D(\boldsymbol{X}_{i,.}, \boldsymbol{Z})} \\ &\leq \sum_{i=1}^{m_1} \min_{j \in m_2} \|\boldsymbol{X}_{i,.} - \boldsymbol{Y}_{j,.}\|_2 + \sum_{i=1}^{m_1} \tilde{D}(\boldsymbol{Y}, \boldsymbol{Z}) \\ &= \tilde{\mathcal{L}}(\boldsymbol{X}, \boldsymbol{Y}) + m_1 \tilde{D}(\boldsymbol{Y}, \boldsymbol{Z}) \end{split}$$

which completes the proof.

Lemma 4 (Javadi and Montanari (2019)). If $A, B \in \mathbb{R}^{m \times n}$, $m \leq n$ are matrices with linearly independent rows, we have

$$\mathcal{L}(A, B)^{1/2} \le 2\kappa(A)D(A, B)^{1/2} + (1 + \sqrt{2})\sqrt{m}D(B, A)^{1/2}$$

where recall that $\kappa(\mathbf{H}) := \sigma_{\max}(\mathbf{H})/\sigma_{\min}(\mathbf{H})$ denotes the condition number of \mathbf{H} .

Lemma 5. Suppose $\boldsymbol{H} \in \mathbb{R}^{k \times n}$, $\boldsymbol{X}_0 \in \mathbb{R}^{m \times n}$ and $\boldsymbol{Z} \in \mathbb{R}^{m \times n}$ is such that $\max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2 \leq \delta$. Let $\boldsymbol{X} = \boldsymbol{X}_0 + \boldsymbol{Z}$. We have

$$\tilde{D}(\boldsymbol{H}, \boldsymbol{X}) \le \tilde{D}(\boldsymbol{H}, \boldsymbol{X}_0) + k\delta \tag{39}$$

$$\tilde{D}(X, H) \le \tilde{D}(X_0, H) + m\delta. \tag{40}$$

Proof. Let us denote the m-dimensional unit simplex by:

$$\Delta^m = \{ \boldsymbol{\alpha} \in \mathbb{R}^m : \sum_{i=1}^m \boldsymbol{\alpha}_i = 1, \boldsymbol{\alpha}_i \geq 0 \}.$$

We have for any $i \in [k]$

$$\begin{split} \tilde{D}(\boldsymbol{H}_{i,.}, \boldsymbol{X}) &= \min_{\boldsymbol{\alpha} \in \Delta^{m}} \|\boldsymbol{H}_{i,.} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{X}_{j,.}\|_{2} \\ &= \min_{\boldsymbol{\alpha} \in \Delta^{m}} \|\boldsymbol{H}_{i,.} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{X}_{j,.}^{0} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{Z}_{j,.}\|_{2} \\ &\leq \min_{\boldsymbol{\alpha} \in \Delta^{m}} \left\{ \|\boldsymbol{H}_{i,.} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{X}_{j,.}^{0}\|_{2} + \|\sum_{j=1}^{m} \alpha_{j} \boldsymbol{Z}_{j,.}\|_{2} \right\} \\ &\leq \min_{\boldsymbol{\alpha} \in \Delta^{m}} \left\{ \|\boldsymbol{H}_{i,.} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{X}_{j,.}^{0}\|_{2} + \sum_{j=1}^{m} \alpha_{j} \|\boldsymbol{Z}_{j,.}\|_{2} \right\} \\ &\leq \min_{\boldsymbol{\alpha} \in \Delta^{m}} \left\{ \|\boldsymbol{H}_{i,.} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{X}_{j,.}^{0}\|_{2} + (\sum_{j=1}^{m} \alpha_{j}) \max_{j \in [m]} \|\boldsymbol{Z}_{j,.}\|_{2} \right\} \\ &\leq \min_{\boldsymbol{\alpha} \in \Delta^{m}} \|\boldsymbol{H}_{i,.} - \sum_{j=1}^{m} \alpha_{j} \boldsymbol{X}_{j,.}^{0}\|_{2} + \delta \\ &= \tilde{D}(\boldsymbol{H}_{i,.}, \boldsymbol{X}_{0}) + \delta. \end{split}$$

Using the above bound, we have:

$$\tilde{D}(\boldsymbol{H}, \boldsymbol{X}) = \sum_{i=1}^{k} \tilde{D}(\boldsymbol{H}_{i,.}, \boldsymbol{X}) \leq \sum_{i=1}^{k} \tilde{D}(\boldsymbol{H}_{i,.}, \boldsymbol{X}_{0}) + k\delta = \tilde{D}(\boldsymbol{H}, \boldsymbol{X}_{0}) + k\delta$$

which establishes (39). We now show (40). For a fixed $i \in [m]$, we have

$$\begin{split} \tilde{D}(\boldsymbol{X}_{i,.}, \boldsymbol{H}) &= \min_{\boldsymbol{\alpha} \in \Delta^k} \|\boldsymbol{X}_{i,.} - \sum_{j=1}^k \boldsymbol{\alpha}_j \boldsymbol{H}_{j,.}\|_2 \\ &= \min_{\boldsymbol{\alpha} \in \Delta^k} \|\boldsymbol{X}_{i,.}^0 - \sum_{j=1}^k \boldsymbol{\alpha}_j \boldsymbol{H}_{j,.} + \boldsymbol{Z}_{i,.}\|_2 \\ &\leq \min_{\boldsymbol{\alpha} \in \Delta^k} \|\boldsymbol{X}_{i,.}^0 - \sum_{j=1}^k \boldsymbol{\alpha}_j \boldsymbol{H}_{j,.}\|_2 + \|\boldsymbol{Z}_{i,.}\|_2 \\ &\leq \min_{\boldsymbol{\alpha} \in \Delta^k} \|\boldsymbol{X}_{i,.}^0 - \sum_{j=1}^k \boldsymbol{\alpha}_j \boldsymbol{H}_{j,.}\|_2 + \delta \\ &= \tilde{D}(\boldsymbol{X}_{i,.}^0, \boldsymbol{H}) + \delta. \end{split}$$

By summing the above bound over i, we arrive at (40).

Lemma 6. Suppose $H \in \mathbb{R}^{k \times n}$, $k \leq n$, has full row rank and $X \in \mathbb{R}^{m \times n}$ is such that $Conv(X) \subseteq Conv(H)$.

Let $\tilde{X} \in \mathbb{R}^{k \times n}$ be a matrix with its *i*-row given by:

$$\tilde{\boldsymbol{X}}_{i,.} = \operatorname*{argmin}_{\boldsymbol{u} \in \{\boldsymbol{X}_{j,.}: j \in [m]\}} \|\boldsymbol{u} - \boldsymbol{H}_{i,.}\|_{2}$$

$$\tag{41}$$

for all $i \in [k]$. Then,

$$\mathcal{L}(\boldsymbol{H}, \boldsymbol{X})^{1/2} \le 2k\kappa(\boldsymbol{H})D(\boldsymbol{H}, \tilde{\boldsymbol{X}})^{1/2}.$$
(42)

Proof. Fix any $\epsilon > 0$. There is a matrix $\mathbf{Z}_{\epsilon} \in \mathbb{R}^{k \times n}$ such that $\|\mathbf{Z}_{\epsilon}\|_{F} \leq \epsilon$ (therefore, $\max_{i \in [m]} \|\mathbf{Z}_{i,.}^{\epsilon}\|_{2} \leq \epsilon$) and $\tilde{\mathbf{X}} + \mathbf{Z}_{\epsilon}$ has full row rank. We have for any $i \in [k]$

$$\tilde{\mathcal{L}}(\boldsymbol{H}_{i,.}, \boldsymbol{X}) = \tilde{\mathcal{L}}(\boldsymbol{H}_{i,.}, \tilde{\boldsymbol{X}})$$

$$= \min_{j \in [k]} \|\boldsymbol{H}_{i,.} - \tilde{\boldsymbol{X}}_{j,.}\|_{2}$$

$$= \min_{j \in [k]} \|\boldsymbol{H}_{i,.} - \tilde{\boldsymbol{X}}_{j,.} + (\boldsymbol{Z}_{\epsilon})_{j,.} - (\boldsymbol{Z}_{\epsilon})_{j,.}\|_{2}$$

$$\leq \min_{j \in [k]} \|\boldsymbol{H}_{i,.} - \tilde{\boldsymbol{X}}_{j,.} - (\boldsymbol{Z}_{\epsilon})_{j,.}\|_{2} + \|(\boldsymbol{Z}_{\epsilon})_{j,.}\|_{2}$$

$$\leq \tilde{\mathcal{L}}(\boldsymbol{H}_{i,.}, \tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon}) + \epsilon.$$
(43)

Note that we have the following inequalities:

$$\mathcal{L}(\boldsymbol{H}, \boldsymbol{X})^{1/2} \stackrel{(a)}{\leq} \tilde{\mathcal{L}}(\boldsymbol{H}, \boldsymbol{X})$$

$$= \sum_{i=1}^{k} \tilde{\mathcal{L}}(\boldsymbol{H}_{i,.}, \boldsymbol{X})$$

$$\stackrel{(b)}{\leq} \tilde{\mathcal{L}}(\boldsymbol{H}, \tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon}) + k\epsilon$$

$$\stackrel{(c)}{\leq} \sqrt{k} \mathcal{L}(\boldsymbol{H}, \tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon})^{1/2} + k\epsilon, \tag{44}$$

where, (a), (c) are results of Lemma 1 and (b) is a result of (43). In addition, by Lemmas 1 and 5, we have:

$$D(\boldsymbol{H}, \tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon})^{1/2} \le \tilde{D}(\boldsymbol{H}, \tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon}) \le \tilde{D}(\boldsymbol{H}, \tilde{\boldsymbol{X}}) + k\epsilon$$
(45)

and

$$D(\tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon}, \boldsymbol{H})^{1/2} \le \tilde{D}(\tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon}, \boldsymbol{H}) \le \tilde{D}(\tilde{\boldsymbol{X}}, \boldsymbol{H}) + m\epsilon = m\epsilon, \tag{46}$$

where the last equality in (46) follows from the fact that $\tilde{D}(\tilde{\boldsymbol{X}}, \boldsymbol{H}) = 0$ (as $\operatorname{Conv}(\boldsymbol{X}) \subseteq \operatorname{Conv}(\boldsymbol{H})$). Starting with (44), we have

$$\mathcal{L}(\boldsymbol{H}, \boldsymbol{X})^{1/2} \le \sqrt{k} \mathcal{L}(\boldsymbol{H}, \tilde{\boldsymbol{X}} + \boldsymbol{Z}_{\epsilon})^{1/2} + k\epsilon \tag{47}$$

$$\leq 2\sqrt{k}\kappa(\boldsymbol{H})D(\boldsymbol{H},\tilde{\boldsymbol{X}}+\boldsymbol{Z}_{\epsilon})^{1/2} + k(1+\sqrt{2})D(\tilde{\boldsymbol{X}}+\boldsymbol{Z}_{\epsilon},\boldsymbol{H})^{1/2} + k\epsilon \tag{48}$$

$$\leq 2\sqrt{k}\kappa(\boldsymbol{H})\tilde{D}(\boldsymbol{H},\tilde{\boldsymbol{X}}) + \epsilon(k + mk(1 + \sqrt{2}) + 2\sqrt{k^3}\kappa(\boldsymbol{H})) \tag{49}$$

$$\leq 2k\kappa(\boldsymbol{H})D(\boldsymbol{H},\tilde{\boldsymbol{X}})^{1/2} + \epsilon(k + mk(1 + \sqrt{2}) + 2\sqrt{k^3}\kappa(\boldsymbol{H}))$$
(50)

where above, inequality (48) is a result of Lemma 4; (49) follows from (45) and (46). Finally, inequality (50) is a result of Lemma 1. As inequality (50) is true for any $\epsilon > 0$, taking the limit $\epsilon \downarrow 0+$, we arrive at (42). \Box

D Proofs of Main Results

D.1 Proof of Theorem 1

Proof. For any i, let us denote:

$$j_i = \operatorname*{argmin}_{j \in [k]} \| \boldsymbol{H}_{i,.} - \boldsymbol{H}_{j,.}^0 \|_2.$$

We have

$$\mathcal{L}(\boldsymbol{H}_{0}, \boldsymbol{H}) = \sum_{i=1}^{k} \min_{j \in [k]} \|\boldsymbol{H}_{i,.}^{0} - \boldsymbol{H}_{j,.}\|_{2}^{2}$$

$$\leq \sum_{i=1}^{k} \|\boldsymbol{H}_{i,.}^{0} - \boldsymbol{H}_{i,.}\|_{2}^{2}$$

$$= \sum_{i=1}^{k} \|\boldsymbol{H}_{i,.}^{0} - \boldsymbol{H}_{j_{i},.}^{0} + \boldsymbol{H}_{j_{i},.}^{0} - \boldsymbol{H}_{i,.}\|_{2}^{2}$$

$$\leq 2 \sum_{i=1}^{k} \|\boldsymbol{H}_{i,.}^{0} - \boldsymbol{H}_{j_{i},.}^{0}\|_{2}^{2} + 2 \sum_{i=1}^{k} \|\boldsymbol{H}_{j_{i},.}^{0} - \boldsymbol{H}_{i,.}\|_{2}^{2}$$

$$\leq 2kb(\boldsymbol{H}_{0})^{2} + 2\mathcal{L}(\boldsymbol{H}, \boldsymbol{H}_{0}),$$

where, the last line follows from the definition of $b(\mathbf{H})$ in (6).

D.2 Proof of Theorem 2

Proof. For any given matrices \boldsymbol{H} and \boldsymbol{H}_0 , there exists a matrix $\boldsymbol{U}_1 \in \{0,1\}^{k \times k}$ such that it has exactly one 1 in each row, such that

$$\mathcal{L}(H, H_0) = \|H - U_1 H_0\|_F^2. \tag{51}$$

In fact, for any $i \in [k]$, we have $U_{i,j}^1 = 1$ where

$$j_i = \underset{j \in [k]}{\operatorname{argmin}} \| \boldsymbol{H}_{i,.} - \boldsymbol{H}_{j,.}^0 \|_2.$$

Noting the noiseless data X_0 is given as $X_0 = W_0 H_0$ where $W_0 \ge 0$ and $W_0 \mathbf{1}_k = \mathbf{1}_m$, one has

$$D(\mathbf{X}_{0}, \mathbf{H})^{1/2} = \min_{\substack{\mathbf{W} \ge 0 \\ \mathbf{W}\mathbf{1}_{k} = \mathbf{1}_{m}}} \|\mathbf{X}_{0} - \mathbf{W}\mathbf{H}\|_{F} \le \|\mathbf{X}_{0} - \mathbf{W}_{0}\mathbf{H}\|_{F}$$

$$= \|\mathbf{W}_{0}\mathbf{H}_{0} - \mathbf{W}_{0}(\mathbf{H} - \mathbf{U}_{1}\mathbf{H}_{0} + \mathbf{U}_{1}\mathbf{H}_{0})\|_{F}$$

$$= \|\mathbf{W}_{0}(\mathbb{I}_{k} - \mathbf{U}_{1})\mathbf{H}_{0} + \mathbf{W}_{0}(\mathbf{H} - \mathbf{U}_{1}\mathbf{H}_{0})\|_{F}$$

$$\le \|\mathbf{W}_{0}(\mathbb{I}_{k} - \mathbf{U}_{1})\mathbf{H}_{0}\|_{F} + \|\mathbf{W}_{0}(\mathbf{H} - \mathbf{U}_{1}\mathbf{H}_{0})\|_{F}.$$
(52)

Note that each row of W_0 sums to 1 and is nonnegative. Therefore, for each row, $\|W_{i,.}^0\|_2^2 \leq 1$ and

$$\|\boldsymbol{W}_0\|_F^2 = \sum_{i=1}^m \|\boldsymbol{W}_{i,.}^0\|_2^2 \le m.$$
 (53)

Using (52), we have the following inequalities:

$$D(\boldsymbol{X}_{0}, \boldsymbol{H})^{1/2} \leq \|\boldsymbol{W}_{0}(\mathbb{I}_{k} - \boldsymbol{U}_{1})\boldsymbol{H}_{0}\|_{F} + \|\boldsymbol{W}_{0}(\boldsymbol{H} - \boldsymbol{U}_{1}\boldsymbol{H}_{0})\|_{F}$$

$$\leq \sqrt{m}\|\boldsymbol{H}_{0}\|_{F}\|\mathbb{I}_{k} - \boldsymbol{U}_{1}\|_{F} + \sqrt{m}\|\boldsymbol{H} - \boldsymbol{U}_{1}\boldsymbol{H}_{0}\|_{F}$$

$$\leq k\sqrt{m}\|\boldsymbol{H}_{0}\|_{F} + \sqrt{m}\mathcal{L}(\boldsymbol{H}, \boldsymbol{H}_{0})^{1/2}$$

where the second inequality makes use of (53); and the last inequality uses (51) and $\|\mathbb{I}_k - U_1\|_F \leq k$ (since, $\mathbb{I}_k - U_1$ is a matrix with every entry between -1 and 1).

Similar to U_1 in the definition (51), we can consider $U_2 \in \{0,1\}^{k \times k}$ such that

$$\mathcal{L}(\boldsymbol{H}_0, \boldsymbol{H}) = \|\boldsymbol{H}_0 - \boldsymbol{U}_2 \boldsymbol{H}\|_F^2. \tag{54}$$

Since $U_2 \mathbf{1}_k = \mathbf{1}_k$; and using the fact that every row of W_0 sums to one, we have: $W_0 U_2 \mathbf{1}_k = \mathbf{1}_m$ and

$$W_0 U_2 \in \{ W : W \ge 0, W \mathbf{1}_k = \mathbf{1}_m \}.$$
 (55)

Note that we have the following:

$$D(\boldsymbol{X}_{0}, \boldsymbol{H})^{1/2} = \min_{\substack{\boldsymbol{W} \geq 0 \\ \boldsymbol{W} \boldsymbol{1}_{k} = \boldsymbol{1}_{m}}} \|\boldsymbol{X}_{0} - \boldsymbol{W}\boldsymbol{H}\|_{F} \overset{(a)}{\leq} \|\boldsymbol{X}_{0} - \boldsymbol{W}_{0}\boldsymbol{U}_{2}\boldsymbol{H}\|_{F}$$
$$= \|\boldsymbol{W}_{0}\boldsymbol{H}_{0} - \boldsymbol{W}_{0}\boldsymbol{U}_{2}\boldsymbol{H}\|_{F}$$
$$\overset{(b)}{\leq} \sqrt{m}\|\boldsymbol{H}_{0} - \boldsymbol{U}_{2}\boldsymbol{H}\|_{F} \overset{(c)}{=} \sqrt{m}\mathcal{L}(\boldsymbol{H}_{0}, \boldsymbol{H})^{1/2},$$

where, (a) uses feasibility condition (55), (b) is due to (53) and (c) uses (54).

D.3 Proof of Theorem 3

Note that the constants in this theorem are listed below.

$$c_{1} = 4\sqrt{k^{3}}\kappa^{2}(\mathbf{H}_{0}) + (1+\sqrt{2})\sqrt{k^{3}}$$

$$c_{2} = 4mk\kappa(\mathbf{H}_{0}) + (1+\sqrt{2})\sqrt{k}(k+\sqrt{k^{3}})$$

$$c_{3} = 2\sqrt{m^{3}}k\kappa(\mathbf{H}_{0}) + (1+\sqrt{2})k^{2}$$

$$c_{4} = k + 2\sqrt{k^{3}}\kappa(\mathbf{H}_{0})$$

$$c_{5} = (k+\sqrt{k^{3}}) + 2mk$$

$$c_{6} = \sqrt{k^{3}} + k\sqrt{m^{3}}$$

$$c_{7} = \frac{\sigma_{\min}(\mathbf{H}_{0})}{6\sqrt{k}}$$

$$c_{8} = 7k\kappa(\mathbf{H}_{0}) + 2(1+\sqrt{2})k^{2}\kappa(\mathbf{H}_{0})$$

$$c_{9} = 7\kappa(\mathbf{H}_{0})(k+\sqrt{k^{3}}) + 2(1+\sqrt{2})\sqrt{k^{3}}m$$

$$c_{10} = 7\kappa(\mathbf{H}_{0})\sqrt{k^{3}} + (1+\sqrt{2})\sqrt{k^{3}}\sqrt{m^{3}}.$$

$$(56)$$

We first present Lemmas 7 and 8 useful for the proof of Theorem 3.

Lemma 7. Under the assumptions of Theorem 3, $P_{\ell}(\mathbf{H}_0)$ is feasible for problem (3).

Proof. By our model-setup, there is a nonnegative W_0 with $W_0 \mathbf{1}_k = \mathbf{1}_m$ such that $X_0 = W_0 H_0$. For $i \in [m]$, let

$$\boldsymbol{v}_i = \boldsymbol{W}_{i..}^0 P_{\ell}(\boldsymbol{H_0}) \in \operatorname{Conv}(P_{\ell}(\boldsymbol{H_0})).$$

Using the definition of v_i above, we have

$$D(X_{0}, P_{\ell}(H_{0})) = \sum_{i=1}^{m} \min_{\mathbf{u} \in \text{Conv}(P_{\ell}(H_{0}))} \|X_{i,.}^{0} - \mathbf{u}\|_{2}^{2} \leq \sum_{i=1}^{m} \|X_{i,.}^{0} - \mathbf{v}_{i}\|_{2}^{2}$$
$$= \|X_{0} - \mathbf{W}_{0}P_{\ell}(H_{0})\|_{F}^{2}.$$
(57)

For $i \in [m]$, we have:

$$D(\mathbf{X}_{i,.}^{0}, P_{\ell}(\mathbf{H}_{0}))^{1/2} \overset{(a)}{\leq} D(\mathbf{X}_{0}, P_{\ell}(\mathbf{H}_{0}))^{1/2}$$

$$\overset{(b)}{\leq} \|\mathbf{X}_{0} - \mathbf{W}_{0} P_{\ell}(\mathbf{H}_{0})\|_{F}$$

$$= \|\mathbf{W}_{0}(\mathbf{H}_{0} - P_{\ell}(\mathbf{H}_{0}))\|_{F}$$

$$= \|\mathbf{W}_{0} P_{\ell}^{\perp}(\mathbf{H}_{0})\|_{F} \leq \|\mathbf{W}_{0}\|_{F} \|P_{\ell}^{\perp}(\mathbf{H}_{0})\|_{F}$$

$$(58)$$

where (a) uses the definition of D(X, H) and (b) is a result of (57). By (53) and (58), we have:

$$D(\mathbf{X}_{i..}^{0}, P_{\ell}(\mathbf{H}_{0}))^{1/2} \le \sqrt{m} \|P_{\ell}^{\perp}(\mathbf{H}_{0})\|_{F} = \beta.$$
 (59)

In what follows, for notational convenience, we denote $\mathbf{H} = P_{\ell}(\mathbf{H}_0)$. Note that \mathbf{H} satisfies the sparsity

constraint in (3). We have the following:

$$D(\boldsymbol{X}_{i,..}, \boldsymbol{H})^{1/2} \leq \tilde{D}(\boldsymbol{X}_{i,..}^{0} + \boldsymbol{Z}_{i,..}, \boldsymbol{H}) \leq \tilde{D}(\boldsymbol{X}_{i,..}^{0}, \boldsymbol{H}) + \|\boldsymbol{Z}_{i,..}\|_{2}$$

$$\leq D(\boldsymbol{X}_{i,..}^{0}, \boldsymbol{H})^{1/2} + \max_{i} \|\boldsymbol{Z}_{i,..}\|_{2}$$

$$\leq \sqrt{m} \|P_{\ell}^{\perp}(\boldsymbol{H}_{0})\|_{F} + \delta = \beta + \delta.$$
(60)

where the first inequality is a result of Lemma 1 and the fact $X = X_0 + Z$; the second inequality is a result of Lemma 5 (we use (40) with m = 1); and the third inequality uses (59). Bound (60) shows that $H = P_{\ell}(H_0)$ is a feasible solution for problem (3).

Lemma 8. Under the assumptions of Theorem 3, one has

$$D(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} \le \sqrt{k} D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)^{1/2} + \sqrt{k^3/m}\beta + (k + \sqrt{k^3})\delta$$
(61)

$$D(\mathbf{H}_{0}, \hat{\mathbf{H}})^{1/2} \leq 2\sqrt{k^{3}}\kappa(\mathbf{H}_{0})D(\mathbf{H}_{0}, \tilde{\mathbf{X}}_{0})^{1/2} + 2mk\delta + mk\beta \cdot (k + \sqrt{k^{3}})\delta.$$
(62)

Proof. We first prove (61). To this end, note that $Conv(X_0) \subseteq Conv(H_0)$, so

$$D(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} \le D(\hat{\boldsymbol{H}}, \boldsymbol{X}_0)^{1/2} \le \tilde{D}(\hat{\boldsymbol{H}}, \boldsymbol{X}_0)$$

where the second inequality is a result of Lemma 1. In addition, as $\hat{\boldsymbol{H}}$ is the optimal solution of (3) and $P_{\ell}(\boldsymbol{H}_0)$ is feasible for problem (3) (by Lemma 7), we have $D(\hat{\boldsymbol{H}}, \boldsymbol{X}) \leq D(P_{\ell}(\boldsymbol{H}_0), \boldsymbol{X})$ and by Lemma 2,

$$\tilde{D}(\hat{\boldsymbol{H}}, \boldsymbol{X}) \le \sqrt{k}\tilde{D}(P_{\ell}(\boldsymbol{H}_0), \boldsymbol{X}). \tag{63}$$

Therefore, one can write

$$D(\hat{\boldsymbol{H}}, \boldsymbol{H}_{0})^{1/2} \overset{(a)}{\leq} \tilde{D}(\hat{\boldsymbol{H}}, \boldsymbol{X}_{0}) = \tilde{D}(\hat{\boldsymbol{H}}, \boldsymbol{X} - \boldsymbol{Z}) \overset{(b)}{\leq} \tilde{D}(\hat{\boldsymbol{H}}, \boldsymbol{X}) + k\delta$$

$$\overset{(c)}{\leq} \sqrt{k} \tilde{D}(P_{\ell}(\boldsymbol{H}_{0}), \boldsymbol{X}) + k\delta = \sqrt{k} \tilde{D}(\boldsymbol{H}_{0} - P_{\ell}^{\perp}(\boldsymbol{H}_{0}), \boldsymbol{X}) + k\delta$$

$$\overset{(d)}{\leq} \sqrt{k} \tilde{D}(\boldsymbol{H}_{0}, \boldsymbol{X}) + \sqrt{k^{3}} \max_{j \in [k]} \|P_{\ell}^{\perp}(\boldsymbol{H}_{0})_{j, \cdot}\|_{2} + k\delta$$

$$\leq \sqrt{k} \tilde{D}(\boldsymbol{H}_{0}, \boldsymbol{X}) + \sqrt{k^{3}} \|P_{\ell}^{\perp}(\boldsymbol{H}_{0})\|_{F} + k\delta$$

$$\leq \sqrt{k} \tilde{D}(\boldsymbol{H}_{0}, \boldsymbol{X}) + \sqrt{k^{3}/m}\beta + k\delta$$

$$\overset{(e)}{\leq} \sqrt{k} \tilde{D}(\boldsymbol{H}_{0}, \boldsymbol{X}_{0}) + \sqrt{k^{3}/m}\beta + (k + \sqrt{k^{3}})\delta$$

$$\overset{(f)}{\leq} \sqrt{k} \tilde{D}(\boldsymbol{H}_{0}, \tilde{\boldsymbol{X}}_{0}) + \sqrt{k^{3}/m}\beta + (k + \sqrt{k^{3}})\delta$$

where (a) is a result of Lemma 1, (b), (d), (e) are results of Lemma 5, (c) is a result of (63), and (f) is true as $\sqrt{k}\tilde{D}(\boldsymbol{H}_0, \boldsymbol{X}_0) \leq \sqrt{k}\tilde{D}(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)$ because $\operatorname{Conv}(\tilde{\boldsymbol{X}}_0) \subseteq \operatorname{Conv}(\boldsymbol{X}_0)$. This establishes (61).

We now proceed to show (62). To this end, using Lemma 3,

$$D(\boldsymbol{H}_0, \hat{\boldsymbol{H}})^{1/2} \le \tilde{D}(\boldsymbol{H}_0, \hat{\boldsymbol{H}}) \le \tilde{\mathcal{L}}(\boldsymbol{H}_0, \boldsymbol{X}_0) + k\tilde{D}(\boldsymbol{X}_0, \hat{\boldsymbol{H}}). \tag{64}$$

Note that by Lemma 6,

$$\tilde{\mathcal{L}}(\boldsymbol{H}_0, \boldsymbol{X}_0) \le 2\sqrt{k^3}\kappa(\boldsymbol{H}_0)D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)^{1/2}.$$
(65)

In addition,

$$\tilde{D}(\boldsymbol{X}_{0}, \hat{\boldsymbol{H}}) = \tilde{D}(\boldsymbol{X} - \boldsymbol{Z}, \hat{\boldsymbol{H}}) \stackrel{(a)}{\leq} \tilde{D}(\boldsymbol{X}, \hat{\boldsymbol{H}}) + m\delta$$

$$\stackrel{(b)}{\leq} \sqrt{m}D(\boldsymbol{X}, \hat{\boldsymbol{H}})^{1/2} + m\delta$$

$$\stackrel{(c)}{\leq} \sqrt{m}\sqrt{\sum_{i=1}^{m} (\delta + \beta)^{2} + m\delta}$$

$$\leq m(\delta + \beta) + m\delta = 2m\delta + m\beta$$
(66)

where (a) is a result of Lemma 5, (b) is a result of Lemma 1 and (c) is due to the constraint $D(X_{i,.}, H)^{1/2} \le \delta + \beta$ in Problem (3). Therefore, by (64), (65) and (66),

$$D(\mathbf{H}_0, \hat{\mathbf{H}})^{1/2} \le 2\sqrt{k^3}\kappa(\mathbf{H}_0)D(\mathbf{H}_0, \tilde{\mathbf{X}}_0)^{1/2} + 2mk\delta + mk\beta$$

which establishes (62). \Box

Proof of Theorem 3. Part 1) If $\hat{\boldsymbol{H}}$ has linearly independent rows, the desired result is achieved by substituting (61) and (62) into Lemma 4 with $\boldsymbol{A} = \boldsymbol{H}_0$ and $\boldsymbol{B} = \hat{\boldsymbol{H}}$. If $\hat{\boldsymbol{H}}$ does not have linearly independent rows, for $\epsilon > 0$ there exists $\hat{\boldsymbol{H}}_{\epsilon}$ with linearly independent rows such that $\|\hat{\boldsymbol{H}}_{\epsilon} - \hat{\boldsymbol{H}}\|_F \leq \epsilon$. Following a path similar to the proof of Lemma 6 and taking the limit $\epsilon \downarrow 0+$ the desired result is achieved. Part 2) By summing (61) and (62),

$$D(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} + D(\boldsymbol{H}_0, \hat{\boldsymbol{H}})^{1/2} \le c_4 D(\boldsymbol{H}_0, \tilde{\boldsymbol{X}}_0)^{1/2} + c_5 \max_{i \in [m]} \|\boldsymbol{Z}_{i,.}\|_2 + c_6 \|P_{\ell}^{\perp}(\boldsymbol{H}_0)\|_F$$

$$\stackrel{(a)}{\le} \frac{\sigma_{\min}(\boldsymbol{H}_0)}{6\sqrt{k}}$$

where (a) is a result of condition (10). As a result,

$$D(\hat{\boldsymbol{H}}, \boldsymbol{H}_0)^{1/2} + D(\boldsymbol{H}_0, \hat{\boldsymbol{H}})^{1/2} \le \frac{\sigma_{\min}(\boldsymbol{H}_0)}{6\sqrt{k}}.$$

Therefore, condition (B.42) of Javadi and Montanari (2019) holds and by Lemma B.3 of Javadi and Montanari (2019), we have

$$\kappa(\hat{\boldsymbol{H}}) \le \frac{7}{2}\kappa(\boldsymbol{H}_0),\tag{67}$$

which shows \hat{H} has linearly independent rows. The rest of the proof is achieved by substituting (67), (61) and (62) into Lemma 4 with $A = \hat{H}$ and $B = H_0$.

D.4 Proof of Proposition 1

The constants in this proposition are listed below:

$$c_{\lambda}^{1} = \left(2\kappa(\boldsymbol{H}_{0})\left[k\sqrt{m}\sqrt{m+\lambda k^{2}} + mk\right] + (1+\sqrt{2})\sqrt{k}\left[\sqrt{mk/\lambda + k^{3}} + k\right]\right)$$

$$c_{\lambda}^{2} = \left(7\kappa(\boldsymbol{H}_{0})\left[\sqrt{mk/\lambda + k^{3}} + k\right] + (1+\sqrt{2})\sqrt{k}\left[k\sqrt{m}\sqrt{m+\lambda k^{2}} + mk\right]\right)$$

$$c_{\lambda}^{3} = \left[(1+m)k + k\sqrt{m}\sqrt{m+\lambda k^{2}} + \sqrt{mk/\lambda + k^{3}}\right].$$
(68)

Lemma 9. Under the assumptions of Proposition 1, one has

$$D(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{H}_{0})^{1/2} \leq \sqrt{k} \sqrt{\frac{m\delta^{2}}{\lambda} + k^{2}\delta^{2}} + k\delta, \tag{69}$$

$$D(\boldsymbol{H}_0, \hat{\boldsymbol{H}}_{\lambda})^{1/2} \le k\sqrt{m}\sqrt{m\delta^2 + \lambda k^2\delta^2} + mk\delta. \tag{70}$$

Proof. Recall that in this proposition, we assume $P_{\ell}(\boldsymbol{H}_0) = \boldsymbol{H}_0$ ($\beta = 0$) and

$$D(\mathbf{H}_0, \tilde{\mathbf{X}}_0) = \tilde{D}(\mathbf{H}_0, \mathbf{X}_0) = 0.$$
(71)

From (60), we have

$$D(\boldsymbol{X}, \boldsymbol{H}_0) = \sum_{i=1}^{m} D(\boldsymbol{X}_{i,.}, \boldsymbol{H}_0) \le m\delta^2.$$
(72)

In addition, we have

$$D(\boldsymbol{H}_0, \boldsymbol{X})^{1/2} \stackrel{(a)}{\leq} \tilde{D}(\boldsymbol{H}_0, \boldsymbol{X}) = \tilde{D}(\boldsymbol{H}_0, \boldsymbol{X}_0 + \boldsymbol{Z}) \stackrel{(b)}{\leq} \tilde{D}(\boldsymbol{H}_0, \boldsymbol{X}_0) + k\delta \stackrel{(c)}{=} k\delta$$
 (73)

where (a) is due to Lemma 1, (b) is due to Lemma 5 and (c) is true because of (71). Therefore, from (73) we have:

$$D(\boldsymbol{H}_0, \boldsymbol{X}) \leq [\tilde{D}(\boldsymbol{H}_0, \boldsymbol{X})]^2 \leq k^2 \delta^2.$$

Let $u = m\delta^2$ and $v = k^2\delta^2$. Note that as $\hat{\boldsymbol{H}}_{\lambda}$ is the optimal solution of the penalized problem (12), by (72) and (73) we have

$$D(X, \hat{H}_{\lambda}) + \lambda D(\hat{H}_{\lambda}, X) \le u + \lambda v. \tag{74}$$

Note that we have the following:

$$D(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{H}_{0})^{1/2} \stackrel{(a)}{\leq} D(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{X}_{0})^{1/2} \stackrel{(b)}{\leq} \tilde{D}(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{X}_{0})$$

$$= \tilde{D}(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{X} - \boldsymbol{Z}) \stackrel{(c)}{\leq} \tilde{D}(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{X}) + k\delta$$

$$\stackrel{(d)}{\leq} \sqrt{k}D(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{X})^{1/2} + k\delta \stackrel{(e)}{\leq} \sqrt{k}\sqrt{\frac{u}{\lambda} + v} + k\delta, \tag{75}$$

where (a) is because $\operatorname{Conv}(\boldsymbol{X}_0) \subseteq \operatorname{Conv}(\boldsymbol{H}_0)$, (b), (d) are due to Lemma 1, (c) is due to Lemma 5; and in (e) we use the observation $D(\hat{\boldsymbol{H}}_{\lambda}, \boldsymbol{X}) \leq u/\lambda + v$ (which follows from (74)). This proves (69).

We will now prove (70). We obtain the following set of inequalities:

$$D(\boldsymbol{H}_{0}, \hat{\boldsymbol{H}}_{\lambda})^{1/2} \overset{(a)}{\leq} \tilde{\mathcal{L}}(\boldsymbol{H}_{0}, \boldsymbol{X}_{0}) + k\tilde{D}(\boldsymbol{X}_{0}, \hat{\boldsymbol{H}})$$

$$\overset{(b)}{\leq} 2\sqrt{k^{3}}\kappa(\boldsymbol{H}_{0})D(\boldsymbol{H}_{0}, \tilde{\boldsymbol{X}}_{0})^{1/2} + k\tilde{D}(\boldsymbol{X}_{0}, \hat{\boldsymbol{H}})$$

$$\overset{(c)}{=} k\tilde{D}(\boldsymbol{X}_{0}, \hat{\boldsymbol{H}}_{\lambda})$$

$$\overset{(d)}{\leq} k\tilde{D}(\boldsymbol{X}, \hat{\boldsymbol{H}}_{\lambda}) + mk\delta$$

$$\overset{(e)}{\leq} k\sqrt{m}D(\boldsymbol{X}, \hat{\boldsymbol{H}}_{\lambda})^{1/2} + mk\delta$$

$$\overset{(f)}{\leq} k\sqrt{m}\sqrt{u + \lambda v} + mk\delta$$

$$(76)$$

where (a) is a result of (64), (b) is due to (65), (c) is due to (71), (d) is due to Lemma 5 with $X_0 = X - Z$, (e) is true by Lemma 1; and (f) is a result of (74). This establishes (70).

Proof of Proposition 1. Part 1) If $\hat{\boldsymbol{H}}_{\lambda}$ has linearly independent rows, this part of the proposition is a direct result of (69) and (70) together with Lemma 4 with $\boldsymbol{A} = \boldsymbol{H}_0$ and $\boldsymbol{B} = \hat{\boldsymbol{H}}_{\lambda}$. If $\hat{\boldsymbol{H}}_{\lambda}$ does not have linearly independent rows, a perturbation argument similar to the proof of Theorem 3 Part 1 suffices. Part 2) Similar to the proof of Theorem 3, condition (14) guarantees $\kappa(\hat{\boldsymbol{H}}_{\lambda}) \leq (7/2)\kappa(\boldsymbol{H}_0)$. The rest of the proof follows from (69) and (70) together with Lemma 4 with $\boldsymbol{A} = \hat{\boldsymbol{H}}_{\lambda}$ and $\boldsymbol{B} = \boldsymbol{H}_0$.

D.5 Proof of Theorem 4

Proof. Part 1) The proof of convergence is based on Theorem 2 of Xu and Yin (2017). Following Xu and Yin (2017), we define the maximum and minimum of Lipschitz constants across three blocks $(\boldsymbol{H}, \boldsymbol{W}, \tilde{\boldsymbol{W}})$ at iteration j as

$$L_{j} = \max\{L_{1}(\boldsymbol{W}_{j}), L_{2}(\boldsymbol{H}_{j}), L_{3}(\boldsymbol{X})\}$$
 and $\ell_{j} = \min\{L_{1}(\boldsymbol{W}_{j}), L_{2}(\boldsymbol{H}_{j}), L_{3}(\boldsymbol{X})\}$

respectively. By substituting the values of $L_1(\mathbf{W}_i), L_2(\mathbf{H}_i), L_3(\mathbf{X}),$

$$L_{j} = \max\{2(\|\boldsymbol{W}_{j}^{T}\boldsymbol{W}_{j}\|_{2} + \lambda), 2\max\{\|\boldsymbol{H}_{j}\boldsymbol{H}_{j}^{T}\|_{2}, \varepsilon\}, 2\lambda\|\boldsymbol{X}\boldsymbol{X}^{T}\|_{2}\}$$
(77)

$$\ell_{j} = \min\{2(\|\boldsymbol{W}_{j}^{T}\boldsymbol{W}_{j}\|_{2} + \lambda), 2\max\{\|\boldsymbol{H}_{j}\boldsymbol{H}_{j}^{T}\|_{2}, \varepsilon\}, 2\lambda\|\boldsymbol{X}\boldsymbol{X}^{T}\|_{2}\}.$$
(78)

As W_j , \tilde{W}_j are simplex matrices and bounded, and considering the cost function $\Psi(.,.,.)$ is bounded from above, H_j needs to be bounded. Consequently, L_j is uniformly bounded from above. In addition, by the assumption $\lambda > 0$, ℓ_j is uniformly bounded away from zero. As a result, the Lipschitz constants across three blocks are uniformly bounded from above and bounded away from zero from below—a condition of Theorem 2 of Xu and Yin (2017). Other conditions of Theorem 2 of Xu and Yin (2017) are satisfied and therefore, this implies the convergence of Algorithm 1.

Part 2) Let

$$T_{i} = \max\{0, H_{i} - [1/L_{1}(W_{i})](-W_{i}^{T}[X - W_{i}H_{i}] + \lambda[H_{i} - \tilde{W}_{i}X])\}.$$
(79)

Note that $T_j \to T$ where T is defined in (20). In addition by the assumption of the second part of the theorem on T, $P_{\ell}(T)$ is unique. First, we show that $P_{\ell}(T_j) \to P_{\ell}(T)$. Let us consider two cases:

1. $\|T\|_0 > \ell$: In this case, there exists $i^* \ge 1$ such that the support of $P_\ell(T_i)$ and $P_\ell(T)$ are the same for $i \ge i^*$. Therefore, for $i \ge i^*$, for $(r, u) \in S(P_\ell(T_i)) = S(P_\ell(T))$ [recall that S(T) is the support of T],

$$P_{\ell}(\boldsymbol{T}_i)_{r,u} = \boldsymbol{T}_{r,u}^i \to \boldsymbol{T}_{r,u} = P_{\ell}(\boldsymbol{T})_{r,u}$$

and for $(r, u) \in S(P_{\ell}(\mathbf{T}_i))^c = S(P_{\ell}(\mathbf{T}))^c$,

$$P_{\ell}(\boldsymbol{T}_i)_{r,u} = 0 = P_{\ell}(\boldsymbol{T})_{r,u}.$$

2. $\|T\|_0 \le \ell$: In this case, $S(T) = S(P_{\ell}(T))$ and there exists $i^* \ge 1$ such that for $i \ge i^*$, $S(T) \subseteq S(P_{\ell}(T_i))$. As a result, for $i \ge i^*$, for $(r, u) \in S(T)$,

$$P_{\ell}(\boldsymbol{T}_i)_{r,u} = \boldsymbol{T}_{r,u}^i \to \boldsymbol{T}_{r,u} = P_{\ell}(\boldsymbol{T})_{r,u}$$

and for $(r, u) \in S(\mathbf{T})^c$,

$$P_{\ell}(T_i)_{r,u} \in \{0, T_{r,u}^i\}$$

and as $T_{r,u}^i \to T_{r,u} = 0$,

$$P_{\ell}(\boldsymbol{T}_i)_{r,u} \to P_{\ell}(\boldsymbol{T})_{r,u} = 0.$$

In addition, note that for any bounded convex set $C \subseteq \mathbb{R}^m$, if $\mathbf{x}_i \to \mathbf{x}^*$, $P_C(\mathbf{x}_i) \to P_C(\mathbf{x}^*)$ where P_C is the projection onto C. This along with the fact that $P_\ell(\mathbf{T}_j) \to P_\ell(\mathbf{T})$, is sufficient to show stationarity:

$$\begin{aligned} \boldsymbol{H}^* &= \lim_{j \to \infty} \boldsymbol{H}_{j+1} \\ &\stackrel{(a)}{=} \lim_{j \to \infty} P_{\ell}(\boldsymbol{T}_{j}) \\ &\stackrel{(b)}{=} P_{\ell}(\boldsymbol{T}) \\ &= P_{\ell}(\max\{0, \boldsymbol{H}^* - [1/L_{1}(\boldsymbol{W}^*)](-\boldsymbol{W}^{*T}[\boldsymbol{X} - \boldsymbol{W}^*\boldsymbol{H}^*] + \lambda[\boldsymbol{H}^* - \tilde{\boldsymbol{W}}^*\boldsymbol{X}])\}) \\ &= \operatorname*{argmin}_{\boldsymbol{H} \geq 0} \left\| \boldsymbol{H} - \left(\boldsymbol{H}^* - \frac{1}{2L_{1}(\boldsymbol{W}^*)} \nabla_{\boldsymbol{H}} \Psi(\boldsymbol{H}^*, \boldsymbol{W}^*, \tilde{\boldsymbol{W}}^*)\right) \right\|_{F}^{2} \end{aligned}$$

where (a) is by the definition of the iterate H_{j+1} and (b) was proved above, showing stationarity for the block H.

D.6 Proof of Proposition 2

Proof. Let $\phi(\boldsymbol{H}, \tilde{\boldsymbol{W}}) = \|\boldsymbol{H} - \tilde{\boldsymbol{W}}\boldsymbol{X}\|_F^2$ be the objective function of (21). Suppose $\boldsymbol{H}^*, \tilde{\boldsymbol{W}}^*$ are optimal solutions to problem (21) and let $\boldsymbol{H}' = \boldsymbol{0}$. Suppose j is such that $\|\boldsymbol{X}_{j,.}\|_2 = \min_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_2$ and $\boldsymbol{e}_j \in \mathbb{R}^m$ is the vector with all coordinates equal to zero except coordinate j equal to one. Let $\tilde{\boldsymbol{W}}' = \mathbf{1}_k \boldsymbol{e}_j^T$. Hence, $\tilde{\boldsymbol{W}}' \boldsymbol{X} = \mathbf{1}_k \boldsymbol{X}_{j,.}$. Note that $\boldsymbol{H}', \tilde{\boldsymbol{W}}'$ are feasible for (21).

We prove the statement of this proposition by the method of contradiction. Suppose there exists $i_0 \in [k]$ such that

$$\|\boldsymbol{H}_{i_{0},.}^{*}\|_{2} > \max_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2} + \sqrt{k} \min_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2}.$$
(80)

Note that for any $v \in \text{Conv}(X)$, there exists $\alpha_1, \dots, \alpha_m \ge 0$ such that they sum to one and $v = \sum_{i=1}^m \alpha_i X_{i,.}$. As a result,

$$\|\boldsymbol{v}\|_{2} = \|\sum_{i=1}^{m} \alpha_{i} \boldsymbol{X}_{i,.}\|_{2} \leq \sum_{i=1}^{m} \alpha_{i} \|\boldsymbol{X}_{i,.}\|_{2} \leq \sum_{i=1}^{m} \alpha_{i} \max_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2} = \max_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2}$$

which when used with (80) leads to:

$$\|\boldsymbol{H}_{i_0,..}^*\|_2 \ge \|\boldsymbol{v}\|_2. \tag{81}$$

One has

$$\begin{split} k \min_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2}^{2} &\stackrel{(a)}{=} \phi(\boldsymbol{H}', \tilde{\boldsymbol{W}}') \\ &\stackrel{(b)}{\geq} \phi(\boldsymbol{H}^{*}, \tilde{\boldsymbol{W}}^{*}) \\ &\geq \|\boldsymbol{H}_{i_{0},.}^{*} - \tilde{\boldsymbol{W}}_{i_{0},.} \boldsymbol{X}\|_{2}^{2} \\ &\geq D(\boldsymbol{H}_{i_{0},.}^{*}, \boldsymbol{X}) = \min_{\boldsymbol{v} \in \operatorname{Conv}(\boldsymbol{X})} \|\boldsymbol{H}_{i_{0},.}^{*} - \boldsymbol{v}\|_{2}^{2} \\ &\stackrel{(c)}{\geq} \min_{\boldsymbol{v} \in \operatorname{Conv}(\boldsymbol{X})} |\|\boldsymbol{H}_{i_{0},.}^{*}\|_{2} - \|\boldsymbol{v}\|_{2}|^{2} \stackrel{(d)}{\geq} |\|\boldsymbol{H}_{i_{0},.}^{*}\|_{2} - \max_{\boldsymbol{v} \in \operatorname{Conv}(\boldsymbol{X})} \|\boldsymbol{v}\|_{2}|^{2} \\ &= (\|\boldsymbol{H}_{i_{0},.}^{*}\|_{2} - \max_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2})^{2} \\ &\stackrel{(e)}{\geq} k \min_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2}^{2}, \end{split}$$

where (a) is true by definition of $\mathbf{H}', \tilde{\mathbf{W}}', (b)$ is due to the optimality of $\mathbf{H}^*, \tilde{\mathbf{W}}^*, (c)$ is true as for any two vectors $\mathbf{a}, \mathbf{b}, \|\mathbf{a} - \mathbf{b}\|_2^2 \ge (\|\mathbf{a}\|_2 - \|\mathbf{b}\|_2)^2$, (d) is due to (81) and (e) is because of (80). This is a contradiction. Hence, for any $i \in [k]$,

$$\|\boldsymbol{H}_{i}^{*}\|_{2} \leq \max_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2} + \sqrt{k} \min_{u \in [m]} \|\boldsymbol{X}_{u,.}\|_{2}.$$

D.7 Proof of Proposition 3

Proof. Part 1) The cost function of (24) is jointly convex in H, \tilde{W}, Z and the feasible set of (24) is convex. Therefore, F is a marginal minimization of a jointly convex function (w.r.t. H, \tilde{W}) over a convex set and is convex (see Section 3.2.5 of Boyd et al. (2004)).

Part 2) For the second part, note that we can rewrite (24) as:

$$F(\boldsymbol{Z}) = \min_{\boldsymbol{H}, \tilde{\boldsymbol{W}}, \boldsymbol{U}} \quad \|\boldsymbol{H} - \boldsymbol{U}\|_F^2$$
s.t. $\boldsymbol{H} \ge 0, \quad \tilde{\boldsymbol{W}} \ge 0, \quad \tilde{\boldsymbol{W}} \boldsymbol{1}_m = \boldsymbol{1}_k$

$$\boldsymbol{H}_{i,j} \le \sqrt{b} \boldsymbol{Z}_{i,j} \quad \forall (i,j) \in [k] \times [n]$$

$$\boldsymbol{U} = \tilde{\boldsymbol{W}} \boldsymbol{X}.$$
(82)

We start by obtaining the dual problem of (82). Note that by enhanced Slater's condition (Boyd et al., 2004), strong duality holds for this problem. The Lagrangian of (82) can be written as

$$L(\boldsymbol{H}, \tilde{\boldsymbol{W}}, \boldsymbol{U}, \boldsymbol{M}, \boldsymbol{\Lambda}, \boldsymbol{\mu}) = \|\boldsymbol{H} - \boldsymbol{U}\|_{F}^{2} - \langle \boldsymbol{M}_{1}, \boldsymbol{H} \rangle - \langle \boldsymbol{M}_{2}, \tilde{\boldsymbol{W}} \rangle + \langle \boldsymbol{M}_{3}, \boldsymbol{U} - \tilde{\boldsymbol{W}} \boldsymbol{X} \rangle$$

$$+ \langle \boldsymbol{\mu}, \tilde{\boldsymbol{W}} \boldsymbol{1}_{m} - \boldsymbol{1}_{k} \rangle + \langle \boldsymbol{\Lambda}, \boldsymbol{H} - \sqrt{b} \boldsymbol{Z} \rangle$$

$$= \left[\|\boldsymbol{H} - \boldsymbol{U}\|_{F}^{2} - \langle \boldsymbol{M}_{1} - \boldsymbol{\Lambda}, \boldsymbol{H} \rangle + \langle \boldsymbol{M}_{3}, \boldsymbol{U} \rangle \right]$$

$$+ \left[-\langle \boldsymbol{M}_{2}, \tilde{\boldsymbol{W}} \rangle - \langle \boldsymbol{M}_{3}, \tilde{\boldsymbol{W}} \boldsymbol{X} \rangle + \langle \boldsymbol{\mu}, \tilde{\boldsymbol{W}} \boldsymbol{1}_{m} \rangle \right]$$

$$+ \left[-\langle \boldsymbol{\mu}, \boldsymbol{1}_{k} \rangle - \langle \boldsymbol{\Lambda}, \sqrt{b} \boldsymbol{Z} \rangle \right],$$
(83)

where $M_1, M_2, M_3, \mu, \Lambda$ are the corresponding Lagrangian variables. By considering the optimality conditions wrt H, U, \tilde{W} , we achieve

$$2(\boldsymbol{H} - \boldsymbol{U}) = \boldsymbol{M}_1 - \boldsymbol{\Lambda},\tag{84}$$

$$2(\boldsymbol{H} - \boldsymbol{U}) = \boldsymbol{M}_3,\tag{85}$$

$$\boldsymbol{M}_2 + \boldsymbol{M}_3 \boldsymbol{X}^T = \boldsymbol{\mu} \boldsymbol{1}_m^T. \tag{86}$$

Using (84), (85), (86) in (83), we get the dual of (82):

$$F(\boldsymbol{Z}) = \max_{\boldsymbol{M}_{1}, \boldsymbol{M}_{2}, \boldsymbol{\Lambda} \geq 0, \boldsymbol{M}_{3}, \boldsymbol{\mu}} - \frac{1}{4} \|\boldsymbol{M}_{3}\|_{F}^{2} - \langle \boldsymbol{\mu}, \mathbf{1}_{k} \rangle - \langle \boldsymbol{\Lambda}, \sqrt{b} \boldsymbol{Z} \rangle$$
s.t. $\boldsymbol{M}_{2} + \boldsymbol{M}_{3} \boldsymbol{X}^{T} = \boldsymbol{\mu} \mathbf{1}_{m}^{T}$

$$\boldsymbol{M}_{1} - \boldsymbol{\Lambda} = \boldsymbol{M}_{3}.$$
(87)

At optimality, note that for $(i, j) \in [k] \times [n]$, $\Lambda_{i,j} = -M_{i,j}^3$ if $M_{i,j}^3 < 0$ and otherwise $\Lambda_{i,j} = 0$ as $\Lambda_{i,j}, Z_{i,j} \ge 0$ and the cost function is higher if $\Lambda_{i,j}$ is smaller. By using Danskin's Theorem (Bertsekas, 1997), if Λ is the optimal solution to (87), $-\sqrt{b}\Lambda$ is a subgradient of F. In addition, at optimality, based on KKT conditions, $M_3 = 2(H - U)$ by (85) and $U = \tilde{W}X$ by feasibility, completing the proof.

References

Vinayak Abrol and Pulkit Sharma. A geometric approach to archetypal analysis via sparse projections. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 42–51. PMLR, 13–18 Jul 2020.

- Sanjeev Arora, Rong Ge, Ravindran Kannan, and Ankur Moitra. Computing a nonnegative matrix factorization—provably. In *Proceedings of the forty-fourth annual ACM symposium on Theory of computing*, pages 145–162, 2012.
- Francis Bach, Julien Mairal, and Jean Ponce. Convex sparse matrix factorizations. arXiv preprint arXiv:0812.1869, 2008.
- Amir Beck and Yonina C. Eldar. Sparsity constrained nonlinear optimization: Optimality conditions and algorithms. SIAM Journal on Optimization, 23(3):1480–1509, 2013.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. SIAM J. Imq. Sci., 2(1):183–202, March 2009.
- Michael W Berry and Murray Browne. Email surveillance using non-negative matrix factorization. Computational & Mathematical Organization Theory, 11(3):249–264, 2005.
- Dimitri P Bertsekas. Nonlinear programming. Journal of the Operational Research Society, 48(3):334–334, 1997.
- Dimitris Bertsimas and Bart Van Parys. Sparse high-dimensional regression: Exact scalable algorithms and phase transitions. *The Annals of Statistics*, 48(1):300–323, 2020.
- Dimitris Bertsimas, Angela King, and Rahul Mazumder. Best subset selection via a modern optimization lens. *The annals of statistics*, pages 813–852, 2016.
- Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- Jean-Philippe Brunet, Pablo Tamayo, Todd R Golub, and Jill P Mesirov. Metagenes and molecular pattern discovery using matrix factorization. *Proceedings of the national academy of sciences*, 101(12):4164–4169, 2004.
- Yuansi Chen, Julien Mairal, and Zaid Harchaoui. Fast and robust archetypal analysis for representation learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1478–1485, 2014.
- M. Chu, F. Diele, R. Plemmons, and S. Ragni. Optimality, computation, and interpretation of nonnegative matrix factorizations. SIAM Journal on Matrix Analysis, pages 4–8030, 2004.
- Adele Cutler and Leo Breiman. Archetypal analysis. Technometrics, 36(4):338–347, 1994.
- Anil Damle and Yuekai Sun. A geometric approach to archetypal analysis and nonnegative matrix factorization. *Technometrics*, 59(3):361–370, 2017.
- David Donoho and Victoria Stodden. When does non-negative matrix factorization give a correct decomposition into parts? In *Advances in Neural Information Processing Systems* 16, pages 1141–1148. MIT Press, 2004.

- John Duchi, Shai Shalev-Shwartz, Yoram Singer, and Tushar Chandra. Efficient projections onto the 11-ball for learning in high dimensions. In *Proceedings of the 25th International Conference on Machine Learning*, ICML '08, page 272–279, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605582054.
- Marco A Duran and Ignacio E Grossmann. An outer-approximation algorithm for a class of mixed-integer nonlinear programs. *Mathematical programming*, 36(3):307–339, 1986.
- E. Elhamifar, G. Sapiro, and R. Vidal. See all by looking at a few: Sparse modeling for finding representative objects. In 2012 IEEE Conference on Computer Vision and Pattern Recognition, pages 1600–1607, 2012.
- Rong Ge and James Zou. Intersecting faces: Non-negative matrix factorization with new guarantees. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2295–2303, Lille, France, 07–09 Jul 2015. PMLR.
- N. Gillis and S. A. Vavasis. Fast and robust recursive algorithmsfor separable nonnegative matrix factorization. IEEE Transactions on Pattern Analysis and Machine Intelligence, 36(4):698–714, 2014. doi: 10.1109/TPAMI.2013.226.
- Edward F. Gonzalez and Yin Zhang. Accelerating the lee-seung algorithm for nonnegative matrix factorization, 2005.
- Hussein Hazimeh and Rahul Mazumder. Fast best subset selection:: Coordinate descent and local combinatorial optimization algorithms. *Operations Research*, 68(5):1517–1537, 2020.
- Hussein Hazimeh, Rahul Mazumder, and Ali Saab. Sparse regression at scale: Branch-and-bound rooted in first-order optimization. arXiv preprint arXiv:2004.06152, 2020.
- Patrik O Hoyer. Non-negative matrix factorization with sparseness constraints. *Journal of machine learning* research, 5(Nov):1457–1469, 2004.
- Peter J Huber. Robust statistics, volume 523. John Wiley & Sons, 2004.
- Hamid Javadi and Andrea Montanari. Nonnegative matrix factorization via archetypal analysis. *Journal of the American Statistical Association*, pages 1–22, 2019.
- Mahdi M Kalayeh, Haroon Idrees, and Mubarak Shah. Nmf-knn: Image annotation using weighted multiview non-negative matrix factorization. In *Proceedings of the IEEE conference on computer vision and* pattern recognition, pages 184–191, 2014.
- Hyunsoo Kim and Haesun Park. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23(12):1495–1502, 05 2007. ISSN 1367-4803.
- Jingu Kim and Haesun Park. Sparse nonnegative matrix factorization for clustering. Technical report, Georgia Institute of Technology, 2008.

- Philip M Kim and Bruce Tidor. Subsystem identification through dimensionality reduction of large-scale gene expression data. *Genome research*, 13(7):1706–1718, 2003.
- Deguang Kong, Chris Ding, and Heng Huang. Robust nonnegative matrix factorization using l21-norm. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 673–682, 2011.
- Dylan Kotliar, Adrian Veres, M Aurel Nagy, Shervin Tabrizi, Eran Hodis, Douglas A Melton, and Pardis C Sabeti. Identifying gene expression programs of cell-type identity and cellular activity with single-cell rna-seq. *Elife*, 8, 2019.
- William H Lawton and Edward A Sylvestre. Self modeling curve resolution. *Technometrics*, 13(3):617–633, 1971.
- Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. Nature, 401(6755):788–791, 1999.
- V. Leplat, A. M. S. Ang, and N. Gillis. Minimum-volume rank-deficient nonnegative matrix factorizations. In ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pages 3402–3406, 2019. doi: 10.1109/ICASSP.2019.8682280.
- C. Lin. Projected gradient methods for nonnegative matrix factorization. *Neural Computation*, 19(10): 2756–2779, 2007.
- Haifeng Liu, Zhaohui Wu, Xuelong Li, Deng Cai, and Thomas S Huang. Constrained nonnegative matrix factorization for image representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(7):1299–1311, 2011.
- Tomohiko Mizutani. Ellipsoidal rounding for nonnegative matrix factorization under noisy separability. *The Journal of Machine Learning Research*, 15(1):1011–1039, 2014.
- Morten Mørup and Lars Kai Hansen. Archetypal analysis for machine learning and data mining. *Neuro-computing*, 80:54 63, 2012. ISSN 0925-2312. Special Issue on Machine Learning for Signal Processing 2010.
- Pentti Paatero and Unto Tapper. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994.
- Robert Peharz and Franz Pernkopf. Sparse nonnegative matrix factorization with l0-constraints. *Neurocomputing*, 80:38–46, 2012.
- Sridhar Ramaswamy, Pablo Tamayo, Ryan Rifkin, Sayan Mukherjee, Chen-Hsiang Yeang, Michael Angelo, Christine Ladd, Michael Reich, Eva Latulippe, Jill P Mesirov, et al. Multiclass cancer diagnosis using tumor gene expression signatures. *Proceedings of the National Academy of Sciences*, 98(26):15149–15154, 2001.
- Ben Recht, Christopher Re, Joel Tropp, and Victor Bittorf. Factoring nonnegative matrices with linear programs. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems* 25, pages 1214–1222. Curran Associates, Inc., 2012.

- F. S. Samaria and A. C. Harter. Parameterisation of a stochastic model for human face identification. In *Proceedings of 1994 IEEE Workshop on Applications of Computer Vision*, pages 138–142, 1994.
- T. Virtanen. Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(3):1066–1074, 2007. doi: 10.1109/TASL.2006.885253.
- Laurence A Wolsey and George L Nemhauser. *Integer and combinatorial optimization*, volume 55. John Wiley & Sons, 1999.
- Jonghye Woo, Jerry L Prince, Maureen Stone, Fangxu Xing, Arnold D Gomez, Jordan R Green, Christopher J Hartnick, Thomas J Brady, Timothy G Reese, Van J Wedeen, et al. A sparse non-negative matrix factorization framework for identifying functional units of tongue behavior from mri. *IEEE transactions* on medical imaging, 38(3):730-740, 2018.
- Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In 2010 IEEE computer society conference on computer vision and pattern recognition, pages 3485–3492. IEEE, 2010.
- Yangyang Xu and Wotao Yin. A globally convergent algorithm for nonconvex optimization based on block coordinate update. *Journal of Scientific Computing*, 72(2):700–734, 2017.
- Z. Yang and E. Oja. Linear and nonlinear projective nonnegative matrix factorization. IEEE Transactions on Neural Networks, 21(5):734-749, 2010. doi: 10.1109/TNN.2010.2041361.