

Online Nonnegative Matrix Factorization with Outliers

Renbo Zhao, *Member, IEEE*, and Vincent Y. F. Tan, *Senior Member, IEEE*

Abstract—We propose a unified and systematic framework for performing online nonnegative matrix factorization in the presence of outliers. Our framework is particularly suited to large-scale data. We propose two solvers based on projected gradient descent and the alternating direction method of multipliers. We prove that the sequence of objective values converges almost surely by appealing to the quasi-martingale convergence theorem. We also show the sequence of learned dictionaries converges to the set of stationary points of the expected loss function almost surely. In addition, we extend our basic problem formulation to various settings with different constraints and regularizers. We also adapt the solvers and analyses to each setting. We perform extensive experiments on both synthetic and real datasets. These experiments demonstrate the computational efficiency and efficacy of our algorithms on tasks such as (parts-based) basis learning, image denoising, shadow removal and foreground-background separation.

Index Terms—Nonnegative matrix factorization, Online learning, Robust learning, Projected gradient descent, Alternating direction method of multipliers

I. INTRODUCTION

In recent years, Nonnegative Matrix Factorization (NMF) has become a popular dimensionality reduction [2] technique, due to its parts-based, non-subtractive interpretation of the learned basis [3]. Given a nonnegative data matrix \mathbf{V} , it seeks to approximately decompose \mathbf{V} into two nonnegative matrices, \mathbf{W} and \mathbf{H} , such that $\mathbf{V} \approx \mathbf{WH}$. In the literature, the fidelity of such an approximation is most commonly measured by $\|\mathbf{V} - \mathbf{WH}\|_F^2$ [4]–[8]. To obtain this approximation, many algorithms have been proposed, including multiplicative updates [4], block principal pivoting [5], projected gradient descent [6], active set method [7], and the alternating direction method of multipliers [8]. These algorithms have promising performances in numerous applications, including document clustering [9], hyperspectral unmixing [10] and audio source separation [11]. However, there are also many studies [12], [13] showing that their performances deteriorate under two common and practical scenarios. The first scenario is when the data matrix \mathbf{V} has a large number of columns (data samples). This situation arises in today’s data-rich environment. *Batch* data processing methods used in the aforementioned algorithms become highly inefficient in terms of the computational time and storage space. The second scenario is the existence of outliers in some of the data samples. For example (e.g.), there are glares and shadows in images due to bad illumination

conditions. Another example is the presence of impulse noises in time series, including speech recordings in natural language processing or temperature recordings in weather forecasting. The outliers, if not handled properly, can significantly corrupt the learned basis matrix, thus the underlying low-dimensional data structure cannot be learned reliably. Moreover, outlier detection and pruning become much more difficult in large-scale datasets [14]. As such, it is imperative to design algorithms that can learn interpretable parts-based basis representations from large-scale datasets whilst being robust to possible outliers.

A. Previous Works

Many efforts have been devoted to address each challenge separately. To handle large datasets, researchers have pursued solutions in three main directions. The first class of algorithms proposed is known as *online NMF* algorithms [12], [15]–[19]. These algorithms aim to refine the basis matrix each time a new data sample is acquired without storing the past data samples. The second class of algorithms is known as *distributed NMF* [20]–[23]. The basic idea behind these algorithms is to distribute the data samples over a network of agents so that several small-scale optimization problems can be performed concurrently. The final class of algorithms is called the *compressed NMF* algorithms [24], [25]. These algorithms perform structured random compression to project the data onto the lower-dimensional manifolds. As such, the size of the dataset can be reduced. These three approaches have successfully reduced the computation and storage complexities—either provably or through numerical experiments. For the existence of outliers, a class of algorithms called the (batch) *robust NMF* [13], [26]–[35] has been proposed to reliably learn the basis matrix by minimizing the effects of the outliers. The robustness against the outliers is achieved via different approaches. These are detailed in Section II-A. However, to the best of our knowledge, so far there are no NMF-based algorithms that are able to systematically handle outliers in large-scale datasets.

B. Main Contributions

In this paper, we propose an algorithm called the *online NMF with outliers* that fills this void. Specifically, our algorithm aims to learn the basis matrix \mathbf{W} in an online manner whilst being robust to outliers. The development of the proposed algorithm involves much more than the straightforward combination or adaptation of online NMF and robust NMF algorithms. Indeed, since there are many ways to “robustify” the NMF algorithms, it is crucial to find an appropriate way

A preliminary work has been published in ICASSP 2016 [1].

The authors are with the Department of Electrical and Computer Engineering and the Department of Mathematics, National University of Singapore. They are supported in part by the NUS Young Investigator Award (grant number R-263-000-B37-133).

to incorporate such robustness guarantees into the online algorithms. Our algorithm proceeds as follows. At each time instant, we solve two optimization problems. The first enables us to learn the coefficient and outlier vectors while the second enables us to update the basis matrix. We propose two solvers based on projected gradient descent (PGD) and alternating direction method of multipliers (ADMM) to solve both optimization problems. Moreover, the presence of outliers also results in more difficulty when we analyze the convergence properties of our algorithms. See Section V-C for a detailed discussion. We remark that in recent years, some algorithms of similar flavors have been proposed, e.g., online robust PCA [36], [37] and online robust dictionary learning [38]. However, due to different problem formulations, our algorithm has many distinctive features, which then calls for different techniques to develop the solvers and analyze the convergence properties. Furthermore, in Section VII, we also observe its superior performance on real-world applications, including (parts-based) basis learning, image denoising, shadow removal and foreground-background separation, over the similar algorithms. In sum, our contributions are threefold:

- 1) We develop two different solvers based on PGD and ADMM to solve the optimization problems at each time instant. These two solvers can be easily extended to two novel solvers for the batch robust NMF problem. The theoretical and empirical performances of both solvers are compared and contrasted.
- 2) Assuming the data are independently drawn from some common distribution \mathbb{P} , we prove the almost sure convergence of the sequence of objective values as well as the almost sure convergence of the sequence of learned basis matrices to the set of stationary points of the expected loss function. The proof techniques involve the use of tools from convex analysis [39] and empirical process theory [40], as well as the quasi-martingale convergence theorem [41, Theorem 9.4 & Proposition 9.5].
- 3) We extend the basic problem setting to various other general settings, by altering the constraint sets and adding regularizers. We also indicate how to adapt our solvers and analyses to each case. By doing so, the applicability of our algorithms is greatly generalized.

C. Paper Organization

This paper is organized as follows. We first provide a more detailed literature survey in Section II. Next we state a formal formulation of our problem in Section III. The algorithms are derived in Section IV and their convergence properties are analyzed in Section V. In Section VI, we extend our basic problem formulations to a wide variety of settings, and indicate how the solvers and analyses can be adapted to each setting. Finally in Section VII, we provide extensive experiment results on both synthetic and real data. The results are compared to those of their batch counterparts and other online matrix factorization algorithms. We conclude the paper in Section VIII stating some promising avenues for further investigations.

In this paper, all the lemmas and sections with indices beginning with ‘S’ will appear in the supplemental material.

D. Notations

In the following, we use capital boldface letters to denote matrices. For example, the (updated) dictionary/basis matrix at time t is denoted by \mathbf{W}_t . We use F and K to denote the ambient dimension and the (known) latent dimension of data respectively. We use lower-case boldface letters to denote vectors. Specifically, at time instant t , we denote the acquired sample vector, learned coefficient vector and outlier vector as \mathbf{v}_t , \mathbf{h}_t and \mathbf{r}_t respectively. For a vector \mathbf{x} , its i -th entry is denoted by x_i . Given a matrix \mathbf{X} , we denote its i -th row as $\mathbf{X}_{i\cdot}$, j -th column as $\mathbf{X}_{\cdot j}$ and (i, j) -th entry by $x_{i,j}$. Moreover, we denote its Frobenius norm by $\|\mathbf{X}\|_F$, spectral norm by $\|\mathbf{X}\|_2$, $\ell_{1,1}$ norm by $\|\mathbf{X}\|_{1,1} \triangleq \sum_{i,j} |x_{i,j}|$ and trace by $\text{tr}(\mathbf{X})$. Inequality $\mathbf{x} \geq 0$ or $\mathbf{X} \geq 0$ denotes entry-wise nonnegativity. We use $\langle \cdot, \cdot \rangle$ to denote the Frobenius inner product between two matrices and $\mathbf{1}$ the vector with all entries equal to one. For a closed convex nonempty set \mathcal{A} , we denote $\mathcal{P}_{\mathcal{A}}$ as the Euclidean projector onto \mathcal{A} . In particular, \mathcal{P}_+ denotes the Euclidean projector onto the nonnegative orthant. Also, the ∞ -indicator function of \mathcal{A} , $I_{\mathcal{A}}$ is defined as

$$I_{\mathcal{A}}(x) \triangleq \begin{cases} 0, & x \in \mathcal{A} \\ \infty, & x \notin \mathcal{A} \end{cases}. \quad (1)$$

In particular, I_+ denotes the ∞ -indicator function of the nonnegative orthant. For $n \in \mathbb{N}$, $[n] := \{1, 2, \dots, n\}$. Also, \mathbb{R}_+ denotes the set of nonnegative real numbers.

II. RELATED WORKS

A. Robust NMF

The canonical NMF problem can be stated as the following minimization problem

$$\min_{\mathbf{W} \in \mathcal{C}, \{\mathbf{h}_i\}_{i=1}^N \geq 0} \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i\|_2^2, \quad (2)$$

where $\mathcal{C} \subseteq \mathbb{R}_+^{F \times K}$ denotes the constraint set for \mathbf{W} and N denotes the number of data samples. In many works [4], [42], [43], \mathcal{C} is set to $\mathbb{R}_+^{F \times K}$. For simplicity, we omit regularizers on \mathbf{W} and $\{\mathbf{h}_i\}_{i=1}^N$ at this point. Since algorithms for the canonical NMF perform unsatisfactorily when the data contain outliers, robust NMF algorithms have been proposed. Previous algorithms for robust NMF fall into two categories. The first category [26]–[33] replaces the (half) squared ℓ_2 loss $1/2 \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i\|_2^2$ in (2) with some other robust loss measure $\psi(\cdot|\cdot)$

$$\min_{\mathbf{W} \in \mathcal{C}, \{\mathbf{h}_i\}_{i=1}^N \geq 0} \frac{1}{N} \sum_{i=1}^N \psi(\mathbf{v}_i | \mathbf{W}\mathbf{h}_i). \quad (3)$$

For example, $\psi(\mathbf{v}_i | \mathbf{W}\mathbf{h}_i)$ can be the ℓ_2 norm $\|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i\|_2$ [31] or the ℓ_1 norm $\|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i\|_1$ [33]. The second category [13], [34], [35] retains the squared ℓ_2 loss but explicitly models the outlier vectors $\{\mathbf{r}_i\}_{i=1}^N$. Specifically, (2) is reformulated as

$$\begin{aligned} \min \quad & \frac{1}{N} \sum_{i=1}^N \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 + \lambda \phi(\mathbf{R}) \\ \text{s. t.} \quad & \mathbf{W} \in \mathcal{C}, \{\mathbf{h}_i\}_{i=1}^N \geq 0, \mathbf{R} \in \mathcal{Q} \end{aligned} \quad (4)$$

where $\lambda \geq 0$ is the regularization parameter, $\mathbf{R} = [\mathbf{r}_1, \dots, \mathbf{r}_N]$ is the outlier matrix, $\phi(\mathbf{R})$ is the regularizer on \mathbf{R} and \mathcal{Q} is the feasible set of \mathbf{R} . Depending on the assumed sparsity structure of \mathbf{R} , $\phi(\mathbf{R})$ can be the $\ell_{2,1}$ norm [34], $\ell_{1,2}$ norm [35] or $\ell_{1,1}$ norm [13] of \mathbf{R} . Robust NMF algorithms typically do not admit strong recovery guarantees of the original data matrix since neither (3) nor (4) are convex programs. However, as shown empirically, the estimated basis matrix $\widehat{\mathbf{W}}$ represents meaningful parts of the data and the residues of the outliers in the reconstructed matrix $\widehat{\mathbf{W}}\widehat{\mathbf{H}}$ are very small [1], [34], [35].

B. Online Matrix Factorization

Existing algorithms on online matrix factorization belong to two distinct categories. The first category of algorithms [16], [18], [36], [37], [44] assumes the data samples $\{\mathbf{v}_t\}_{t \geq 1}$ are generated independently from a time-invariant distribution \mathbb{P} . Under this assumption, it is possible to provide theoretical guarantees on the convergence of the online stochastic algorithms by leveraging the empirical process theory, as was done in [16], [36], [37], [44]. These methods have extensive applications, including document clustering [18], image inpainting [44], face recognition [16], and image annotation [16]. The second category of algorithms [12], [15], [17], [19], [38], [45]–[48], with major applications in visual tracking, assumes that the data generation distribution \mathbb{P} is time-varying. Although these assumptions are weaker than those in the first class of algorithms, it is very difficult to provide theoretical guarantees on the convergence of the online algorithms.

C. Online Low-rank and Sparse Modeling

Another related line of works [49]–[53] aims to recover the low-rank data matrix \mathbf{L} and the sparse outlier matrix \mathbf{S} from their additive mixture \mathbf{M} in an online fashion. Among these works, [49]–[51] assume that the sequence of mixture vectors (columns of \mathbf{M}) arrives in a streaming fashion. The authors derive the recovery algorithms based on their proposed models of the sequence of ground-truth data vectors (columns of \mathbf{L}) and outlier vectors (columns of \mathbf{S}). The authors of [52], [53] adopt a different approach. They allow the number of mixture vectors to be large but finite and use an alternating minimization approach to solve a variant of (4) (with additional Tikhonov regularizers on \mathbf{W} and $\{\mathbf{h}_i\}_{i=1}^N$). In learning the coefficient vectors $\{\mathbf{h}_i\}_{i=1}^N$ (with fixed \mathbf{W}), the authors employ the stochastic gradient descent method, thus ensuring that their algorithms are scalable to large-scale data.

III. PROBLEM FORMULATION

Following Section II-A, in this work we explicitly model the outlier vectors as $\{\mathbf{r}_t\}_{t \geq 1}$. Also, we assume the data generation distribution \mathbb{P} is time-invariant, for ease of the convergence analysis. First, for a fixed data sample \mathbf{v} and a fixed basis matrix \mathbf{W} , define the loss function with respect to (w.r.t.) \mathbf{v} and \mathbf{W} , $\ell(\mathbf{v}, \mathbf{W})$ as

$$\ell(\mathbf{v}, \mathbf{W}) \triangleq \min_{\mathbf{h} \geq 0, \mathbf{r} \in \mathcal{R}} \tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r}), \quad (5)$$

where

$$\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r}) \triangleq \frac{1}{2} \|\mathbf{v} - \mathbf{W}\mathbf{h} - \mathbf{r}\|_2^2 + \lambda \|\mathbf{r}\|_1, \quad (6)$$

and $\mathcal{R} \triangleq \{\mathbf{r} \in \mathbb{R}^F \mid \|\mathbf{r}\|_\infty \leq M\}$ is the constraint set of the outlier vector \mathbf{r} . Here we use the ℓ_1 regularization to promote entrywise sparsity on \mathbf{r} .

Next, given a finite set of data samples $\{\mathbf{v}_i\}_{i \in [t]} \stackrel{iid}{\sim} \mathbb{P}$, we define the *empirical loss* associated with $\{\mathbf{v}_i\}_{i \in [t]}$, $f_t(\mathbf{W})$ as

$$f_t(\mathbf{W}) \triangleq \frac{1}{t} \sum_{i=1}^t \ell(\mathbf{v}_i, \mathbf{W}) \quad (7)$$

where $\mathcal{C} \triangleq \{\mathbf{W} \in \mathbb{R}_+^{F \times K} \mid \|\mathbf{W}_{:,i}\|_2 \leq 1, \forall i \in [K]\}$.

Following the convention of the online learning literature [54]–[56], instead of minimizing the empirical loss $f_t(\mathbf{W})$ in (7), we aim to minimize the *expected loss* $f(\mathbf{W})$, i.e.,

$$\min_{\mathbf{W} \in \mathcal{C}} \left[f(\mathbf{W}) \triangleq \mathbb{E}_{\mathbf{v} \sim \mathbb{P}} [\ell(\mathbf{v}, \mathbf{W})] \right]. \quad (8)$$

In other words, we aim to solve a (non-convex) stochastic program [57]. Note that by the strong law of large numbers, given any $\mathbf{W} \in \mathcal{C}$, we have $f_t(\mathbf{W}) \xrightarrow{a.s.} f(\mathbf{W})$ as $t \rightarrow \infty$.

Remark 1. We make three remarks here. First we explain the reasonings behind the choice of the constraint sets \mathcal{C} and \mathcal{R} . The set \mathcal{C} constrains the columns of \mathbf{W} in the unit (nonnegative) ℓ_2 ball. This is to prevent the entries of \mathbf{W} from being unbounded, following the conventions in [6], [58]. The set \mathcal{R} uniformly bounds the entries of \mathbf{r} . This is because in practice, both the underlying data (without outliers) and the observed data are uniformly bounded entrywise. Since we do not require \mathbf{r} to be exactly recovered, this prior information can often improve the estimation of \mathbf{r} . For real data, the bound $M > 0$ can often be easily chosen. For example, for gray-scale images, M can be chosen as $2^m - 1$, where m is the number of bits per pixel. In the case of matrices containing ratings from users with a maximum rating of v , M can be chosen as v . In some scenarios where M is difficult to estimate, we simply set $M = \infty$. As will be shown in Sections IV and V, the algorithms and analyses developed for finite M can be easily adapted to infinite M . Second, for the sake of brevity, we omit regularizing \mathbf{W} and \mathbf{h} . Such regularizations, together with other possible constraint sets of \mathbf{W} and \mathbf{r} will be discussed in Section VI. Third, we assume the ambient data dimension F , the latent data dimension K and the penalty parameter λ in (5) are time-invariant parameters.¹

IV. ALGORITHMS

To tackle the problem proposed in Section III, we leverage the *stochastic majorization-minimization (MM)* framework [60], [61], which has been widely used in previous works on online matrix factorization [12], [15], [17], [19], [36]–[38], [44], [45]. In essence, such framework decomposes the optimization problem in (8) into two steps, namely *nonnegative encoding* and *dictionary update*. Concretely, at a time instant t

¹In this work, we do not simultaneously consider the data with high ambient dimensions. An attempt on this problem in the context of dictionary learning with the squared- ℓ_2 loss has been made in [59].

($t \geq 1$), we first learn the coefficient vector \mathbf{h}_t and the outlier vector \mathbf{r}_t based on the newly acquired data sample \mathbf{v}_t and the previous dictionary matrix \mathbf{W}_{t-1} . Specifically, we solve the following convex optimization problem

$$(\mathbf{h}_t, \mathbf{r}_t) = \arg \min_{\mathbf{h} \geq 0, \mathbf{r} \in \mathcal{R}} \tilde{\ell}(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}, \mathbf{r}). \quad (9)$$

Here the initial basis matrix \mathbf{W}_0 is randomly chosen in \mathcal{C} . Next, based on the past statistics $\{\mathbf{v}_i, \mathbf{h}_i, \mathbf{r}_i\}_{i \in [t]}$, the basis matrix is updated to

$$\mathbf{W}_t = \arg \min_{\mathbf{W} \in \mathcal{C}} \tilde{f}_t(\mathbf{W}), \quad (10)$$

where

$$\tilde{f}_t(\mathbf{W}) \triangleq \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 + \lambda \|\mathbf{r}_i\|_1. \quad (11)$$

We note that (10) can be rewritten as

$$\mathbf{W}_t = \arg \min_{\mathbf{W} \in \mathcal{C}} \frac{1}{2} \text{tr}(\mathbf{W}^T \mathbf{W} \mathbf{A}_t) - \text{tr}(\mathbf{W}^T \mathbf{B}_t), \quad (12)$$

where $\mathbf{A}_t \triangleq 1/t \sum_{i=1}^t \mathbf{h}_i \mathbf{h}_i^T$ and $\mathbf{B}_t \triangleq 1/t \sum_{i=1}^t (\mathbf{v}_i - \mathbf{r}_i) \mathbf{h}_i^T$ are the sufficient statistics. From (12), we observe that our algorithm has a storage complexity independent of t since only \mathbf{A}_t , \mathbf{B}_t and \mathbf{W}_t need to be stored and updated.

To solve (9) and (12) (at a fixed time instant t), we propose two solvers based on PGD and ADMM respectively. For ease of reference, we refer to the former algorithm as OPGD and the latter as OADMM. We now explain the motivations behind proposing these two solvers. Since both (9) and (12) are constrained optimization problems, the most straightforward solver would be based on PGD. Although such a solver has a linear computational complexity per iteration, it typically needs a large number of iterations to converge. Moreover, it is also easily trapped in bad local minima [62]. Thus, it would be meaningful to contrast its performance with a solver with very different properties. This leads us to propose another solver based on ADMM. Such a solver has a higher computational complexity per iteration but typically needs fewer number of iterations to converge [8]. It is also less susceptible to bad local minima since it solves optimization problems in the dual space. In Section VII, we will show that the practical performances of these two solvers are comparable despite the different properties they possess. Thus either solver can be used for most practical purposes.

In the sequel, we omit the time subscript t to keep notations uncluttered. In the iterations, the updated value of a variable is denoted with the superscript '+'. Pseudo-codes of the entire algorithm (for N data samples) are provided in Algorithm 1.

A. Online Algorithm Based on PGD (OPGD)

1) *PGD solver for (9)*: For a fixed \mathbf{W} , we solve (9) by alternating between the following two steps

$$\mathbf{h}^+ := \arg \min_{\mathbf{h} \geq 0} Q_\eta(\mathbf{h}' | \mathbf{h}), \quad (13)$$

$$\mathbf{r}^+ := \arg \min_{\mathbf{r}' \in \mathcal{R}} \frac{1}{2} \|\mathbf{v} - \mathbf{W}\mathbf{h}^+ - \mathbf{r}'\|_2^2 + \lambda \|\mathbf{r}'\|_1, \quad (14)$$

where²

$$Q_\eta(\mathbf{h}' | \mathbf{h}) \triangleq q(\mathbf{h}) + \langle \nabla q(\mathbf{h}), \mathbf{h}' - \mathbf{h} \rangle + \frac{1}{2\eta} \|\mathbf{h}' - \mathbf{h}\|_2^2,$$

$q(\mathbf{h}) \triangleq \frac{1}{2} \|\mathbf{v} - \mathbf{W}\mathbf{h} - \mathbf{r}\|_2^2$ and $\eta \in (0, 1/L]$ with $L \triangleq \|\mathbf{W}\|_2^2$. The steps (13) and (14) can be interpreted based on the framework of *block MM* [63], [64]. Specifically, it is easy to verify that the steps (13) and (14) amount to finding the (unique) minimizers of the majorant functions³ of $\mathbf{h}' \mapsto \ell(\mathbf{v}, \mathbf{W}, \mathbf{h}', \mathbf{r})$ at \mathbf{h} and $\mathbf{r}' \mapsto \tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}^+, \mathbf{r}')$ at \mathbf{r} respectively. Therefore, the convergence analysis in [63] guarantees that such alternating minimization procedure converges to a global optimum of (9).

In addition, we notice that the minimizations in both (13) and (14) have closed-form solutions. For (13), the solution is given by the PGD update step (with constant step size)

$$\mathbf{h}^+ := \mathcal{P}_+(\mathbf{h} - \eta \nabla q(\mathbf{h})). \quad (15)$$

For ease of parameter tuning, we rewrite $\eta = \bar{\kappa}/L$ ($0 < \bar{\kappa} \leq 1$). Furthermore, we fix η (or $\bar{\kappa}$) throughout all iterations.

For (14), if $M = \infty$, the solution is precisely given by

$$\mathbf{r}^+ := \mathcal{S}_\lambda(\mathbf{v} - \mathbf{W}\mathbf{h}^+), \quad (16)$$

where \mathcal{S}_λ is the (elementwise) soft-thresholding operator threshold λ . Otherwise, when M is finite, using [65, Lemma 5] (see Lemma S-3), we have

$$\mathbf{r}^+ := \tilde{\mathcal{S}}_{\lambda, M}(\mathbf{v} - \mathbf{W}\mathbf{h}^+), \quad (17)$$

where for any $\mathbf{x} \in \mathbb{R}^F$ and $i \in [F]$,

$$\left(\tilde{\mathcal{S}}_{\lambda, M}(\mathbf{x}) \right)_i := \begin{cases} 0, & |x_i| < \lambda \\ x_i - \text{sgn}(x_i)\lambda, & \lambda \leq |x_i| \leq \lambda + M \\ \text{sgn}(x_i)M, & |x_i| > \lambda + M \end{cases}.$$

2) *PGD solver for (12)*: Similar to the procedure for solving (13), we first rewrite (12) as

$$\min_{\mathbf{W}} p_t(\mathbf{W}) + I_C(\mathbf{W}), \text{ where}$$

$$p_t(\mathbf{W}) = \frac{1}{2} \text{tr}(\mathbf{W}^T \mathbf{W} \mathbf{A}_t) - \text{tr}(\mathbf{W}^T \mathbf{B}_t). \quad (18)$$

First, it is easy to see p_t is convex and differentiable, ∇p_t is Lipschitz with constant $\tilde{L}_t \triangleq \|\mathbf{A}_t\|_F$. Thus, we can construct a majorant function $P_t(\mathbf{W}' | \mathbf{W})$ for $p_t(\mathbf{W}')$ at $\mathbf{W} \in \mathcal{C}$

$$P_t(\mathbf{W}' | \mathbf{W}) = p_t(\mathbf{W}) + \langle \nabla p_t(\mathbf{W}), \mathbf{W}' - \mathbf{W} \rangle + \frac{1}{2\tilde{\eta}_t} \|\mathbf{W}' - \mathbf{W}\|_F^2, \quad (19)$$

where $\nabla p_t(\mathbf{W}) = \mathbf{W}\mathbf{A}_t - \mathbf{B}_t$ and $\tilde{\eta}_t \in (0, 1/\tilde{L}_t]$. Minimizing $P_t(\mathbf{W}' | \mathbf{W}) + I_C(\mathbf{W}')$ over $\mathbf{W}' \in \mathcal{C}$, we have

$$\mathbf{W}^+ := \arg \min_{\mathbf{W}' \in \mathcal{C}} \|\mathbf{W}' - (\mathbf{W} - \tilde{\eta}_t \nabla p_t(\mathbf{W}))\|_F^2 \quad (20)$$

$$:= \mathcal{P}_C(\mathbf{W} - \tilde{\eta}_t \nabla p_t(\mathbf{W})). \quad (21)$$

²At time t , the value of L is computed based on \mathbf{W}_{t-1} , i.e., $L_t \triangleq \|\mathbf{W}_{t-1}\|_2^2$.

³For a function g with domain \mathcal{G} , its majorant at $\kappa \in \mathcal{G}$, \tilde{g} is the function that satisfies i) $\tilde{g} \geq g$ on \mathcal{G} and ii) $\tilde{g}(\kappa) = g(\kappa)$.

By Lemma S-4, the projection step in (21) is given by

$$\mathbf{W}_{:j}^+ := \frac{\mathcal{P}_+(\mathbf{W} - \tilde{\eta}_t \nabla p_t(\mathbf{W}))_{:j}}{\max\{1, \|\mathcal{P}_+(\mathbf{W} - \tilde{\eta}_t \nabla p_t(\mathbf{W}))_{:j}\|_2\}}, \quad \forall j \in [K]. \quad (22)$$

Again, for each iteration given in (21), we use the same step size $\tilde{\eta}_t = \tilde{\kappa}_t/L$ where $0 < \tilde{\kappa}_t \leq 1$.

Remark 2. In the literature [66], [67], the accelerated proximal gradient descent (APGD) method has been proposed to accelerate the canonical PGD method. With an additional extrapolation step in each iteration, the convergence rate can be improved from $O(1/k)$ to $O(1/k^2)$, where k denotes the number of iterations. However, since in our implementations, both (9) and (12) were only solved to a prescribed accuracy (see Section VII-B2), we observed no significant reduction in running times on the real tasks. Thus for simplicity, we only employ the canonical PGD method.

B. Online Algorithm Based on ADMM (OADMM)

Below we only present the update rules derived via ADMM. The detailed derivation steps are shown in Section S-1.

1) *ADMM solver for (9):* First we reformulate (9) as

$$\begin{aligned} \min_{\mathbf{h}, \mathbf{r}} \quad & \frac{1}{2} \|\mathbf{v} - \mathbf{W}\mathbf{h} - \mathbf{r}\|_2^2 + \lambda \|\mathbf{r}\|_1 + I_+(\mathbf{u}) + I_{\mathcal{R}}(\mathbf{q}) \\ \text{s. t.} \quad & \mathbf{h} = \mathbf{u}, \mathbf{r} = \mathbf{q} \end{aligned} \quad (23)$$

Thus the augmented Lagrangian is

$$\begin{aligned} \mathcal{L}(\mathbf{h}, \mathbf{r}, \mathbf{u}, \mathbf{q}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = & \frac{1}{2} \|\mathbf{v} - \mathbf{W}\mathbf{h} - \mathbf{r}\|_2^2 + \lambda \|\mathbf{r}\|_1 \\ & + I_+(\mathbf{u}) + I_{\mathcal{R}}(\mathbf{q}) + \boldsymbol{\alpha}^T (\mathbf{h} - \mathbf{u}) + \boldsymbol{\beta}^T (\mathbf{r} - \mathbf{q}) \\ & + \frac{\rho_1}{2} \|\mathbf{h} - \mathbf{u}\|_2^2 + \frac{\rho_2}{2} \|\mathbf{r} - \mathbf{q}\|_2^2, \end{aligned} \quad (24)$$

where $\boldsymbol{\alpha}, \boldsymbol{\beta}$ are dual variables and ρ_1, ρ_2 are positive penalty parameters. Then we sequentially update each of $\mathbf{h}, \mathbf{r}, \mathbf{u}, \mathbf{q}, \boldsymbol{\alpha}$ and $\boldsymbol{\beta}$ while keeping other variables fixed

$$\mathbf{h}^+ := (\mathbf{W}^T \mathbf{W} + \rho_1 \mathbf{I})^{-1} (\mathbf{W}^T (\mathbf{v} - \mathbf{r}) + \rho_1 \mathbf{u} - \boldsymbol{\alpha}) \quad (25)$$

$$\mathbf{r}^+ := \mathcal{S}_\lambda(\rho_2 \mathbf{q} + \mathbf{v} - \boldsymbol{\beta} - \mathbf{W}\mathbf{h}) / (1 + \rho_2) \quad (26)$$

$$\mathbf{u}^+ := \mathcal{P}_+(\mathbf{h}^+ + \boldsymbol{\alpha} / \rho_1) \quad (27)$$

$$\mathbf{q}^+ := \mathcal{P}_{\mathcal{R}}(\mathbf{r}^+ + \boldsymbol{\beta} / \rho_2) \quad (28)$$

$$\boldsymbol{\alpha}^+ := \boldsymbol{\alpha} + \rho_1 (\mathbf{h}^+ - \mathbf{u}^+) \quad (29)$$

$$\boldsymbol{\beta}^+ := \boldsymbol{\beta} + \rho_2 (\mathbf{r}^+ - \mathbf{q}^+). \quad (30)$$

2) *ADMM solver for (12):* Again, we rewrite (12) as

$$\begin{aligned} \min \quad & \frac{1}{2} \text{tr}(\mathbf{W}^T \mathbf{W} \mathbf{A}_t) - \text{tr}(\mathbf{W}^T \mathbf{B}_t) + I_C(\mathbf{Q}) \\ \text{s. t.} \quad & \mathbf{W} = \mathbf{Q} \end{aligned} \quad (31)$$

The augmented Lagrangian in this case is

$$\begin{aligned} \mathcal{L}(\mathbf{W}, \mathbf{Q}, \mathbf{D}) = & \frac{1}{2} \text{tr}(\mathbf{W}^T \mathbf{W} \mathbf{A}_t) - \text{tr}(\mathbf{W}^T \mathbf{B}_t) + I_C(\mathbf{Q}) \\ & + \langle \mathbf{D}, \mathbf{W} - \mathbf{Q} \rangle + \frac{\rho_3}{2} \|\mathbf{W} - \mathbf{Q}\|_F^2, \end{aligned} \quad (32)$$

where \mathbf{D} is the dual variable and ρ_3 is a positive penalty parameter. Minimizing \mathbf{W}, \mathbf{Q} and \mathbf{D} sequentially yields

$$\mathbf{W}^+ := (\mathbf{B}_t - \mathbf{D} + \rho_3 \mathbf{Q}) (\mathbf{A}_t + \rho_3 \mathbf{I})^{-1} \quad (33)$$

Algorithm 1 Online NMF with outliers (ONMFO)

Input: Data samples $\{\mathbf{v}_i\}_{i \in [N]}$, penalty parameter λ , initial dictionary matrix \mathbf{W}_0

Initialize sufficient statistics: $\mathbf{A}_0 := \mathbf{0}, \mathbf{B}_0 := \mathbf{0}$

for $t = 1$ to N **do**

1) Acquire a data sample \mathbf{v}_t .

2) Learn the coefficient vector \mathbf{h}_t and the outlier vector \mathbf{r}_t based on \mathbf{W}_{t-1} , using the solvers based on PGD or ADMM (detailed in Sections IV-A1 and IV-B1)

$$(\mathbf{h}_t, \mathbf{r}_t) := \arg \min_{\mathbf{h} \geq \mathbf{0}, \mathbf{r} \in \mathcal{R}} \tilde{\ell}(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}, \mathbf{r}). \quad (36)$$

3) Update the sufficient statistics

$$\mathbf{A}_t := 1/t \{(t-1)\mathbf{A}_{t-1} + \mathbf{h}_t \mathbf{h}_t^T\},$$

$$\mathbf{B}_t := 1/t \{(t-1)\mathbf{B}_{t-1} + (\mathbf{v}_t - \mathbf{r}_t) \mathbf{h}_t^T\}.$$

4) Learn the dictionary matrix \mathbf{W}_t based on \mathbf{A}_t and \mathbf{B}_t , using the solvers based on PGD or ADMM (detailed in Section IV-A2 and IV-B2)

$$\mathbf{W}_t := \arg \min_{\mathbf{W} \in \mathcal{C}} \frac{1}{2} \text{tr}(\mathbf{W}^T \mathbf{W} \mathbf{A}_t) - \text{tr}(\mathbf{W}^T \mathbf{B}_t) \quad (37)$$

end for

Output: Final dictionary matrix \mathbf{W}_N

$$\mathbf{Q}^+ := \mathcal{P}_C(\mathbf{W}^+ + \mathbf{D} / \rho_3) \quad (34)$$

$$\mathbf{D}^+ := \mathbf{D} + \rho_3 (\mathbf{W}^+ - \mathbf{Q}^+). \quad (35)$$

Remark 3. A few comments are now in order: First, although algorithms based on both PGD and ADMM exist in the literature for the canonical NMF problem (see for example, [6], [8], [62]), our problem setting is different from those in the previous works. Specifically, our problem explicitly models the outlier vector \mathbf{r} and the constraint set on \mathbf{W} is more complicated than the non-negative orthant $\mathbb{R}_+^{F \times K}$. Moreover, for the algorithm based on PGD, we use a fixed step size though all the iterations. In contrast, the usual projected gradient method applied to NMF [6] involves computationally intensive Armijo-type line searches for the step sizes. Second, although updates (25) and (33) involve matrix inversions, the operations do not incur high computational complexity since both matrices to be inverted have sizes $K \times K$ where $K \ll F$. Third, the proposed two solvers can be easily extended to solve the batch NMF problem with outliers. Thus, we have also effectively proposed *two new solvers* for the (batch) robust NMF problem. These extensions are detailed in Section S-2. In the sequel we term these two batch algorithms as BPGD and BADMM respectively. Finally, the stopping criteria of both OPGD and OADMM for solving (9) and (12) will be described in Section VII-B2.

V. CONVERGENCE ANALYSES

In this section, we analyze the convergence of both the sequence of objective values $\{f(\mathbf{W}_t)\}_{t \geq 1}$ (Theorem 1) and the sequence of dictionary matrices $\{\mathbf{W}_t\}_{t \geq 1}$ (Theorem 2) produced by Algorithm 1. It is worth noticing that the analyses

are independent of the specific solvers used in steps 2) and 4) in Algorithm 1. Indeed, in our analyses, we only leverage the fact that both (9) and (12) can be exactly solved.

A. Preliminaries

First, given any $\mathbf{h}' \geq 0$ and $\mathbf{r}' \in \mathcal{R}$, we observe that $(\mathbf{v}, \mathbf{W}) \mapsto \tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}', \mathbf{r}')$ and \tilde{f}_t serve as upper-bound functions for ℓ (on $\mathbb{R}_+^F \times \mathcal{C}$) and f_t (on \mathcal{C}) respectively. Denote $(\mathbf{h}^*, \mathbf{r}^*)$ as an optimal solution of (5). (Note that $(\mathbf{h}^*, \mathbf{r}^*)$ is a function of \mathbf{v} and \mathbf{W} and always exists since $\tilde{\ell}(\mathbf{v}', \mathbf{W}', \mathbf{h}, \mathbf{r})$ is closed and coercive on $\mathbb{R}_+^K \times \mathcal{R}$ for any $\mathbf{v}' \in \mathbb{R}_+^F$ and $\mathbf{W}' \in \mathcal{C}$.) Clearly we have $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}^*, \mathbf{r}^*) = \ell(\mathbf{v}, \mathbf{W})$, for any $(\mathbf{v}, \mathbf{W}) \in \mathbb{R}_+^F \times \mathcal{C}$.

Next, we notice that ℓ can be equivalently defined as

$$\ell(\mathbf{v}, \mathbf{W}) = \min_{(\mathbf{h}, \mathbf{r}) \in \mathcal{H} \times \mathcal{R}} \tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r}), \quad (38)$$

where \mathcal{H} is a compact and convex set in \mathbb{R}_+^K . This is because \mathbf{h}^* in (5) is bounded, due to the boundedness of \mathbf{v}_i and \mathbf{W} . Thus it suffices to consider minimizing \mathbf{h} over a compact set in \mathbb{R}_+^K . We can further choose \mathcal{H} to be convex. If $M = \infty$, we can similarly show \mathbf{r}^* is bounded thus still consider minimizing \mathbf{r} over a compact and convex set \mathcal{R} in \mathbb{R}^F .

We make three assumptions for our subsequent analyses.

Assumptions.

- 1) The data generation distribution \mathbb{P} has a compact support set \mathcal{V} .
- 2) For all $(\mathbf{v}, \mathbf{W}) \in \mathcal{V} \times \mathcal{C}$, $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$ is jointly m_1 -strongly convex in (\mathbf{h}, \mathbf{r}) for some constant $m_1 > 0$.
- 3) For all $t \geq 1$, $\tilde{f}_t(\mathbf{W})$ is m_2 -strongly convex on \mathcal{C} for some constant $m_2 > 0$.

Remark 4. All of the assumptions above are made to simplify the convergence analyses. Among them, Assumption 1 naturally holds for real data, which are always bounded. Assumptions 2 and 3 play roles in proving Lemma 1 and 2 respectively. Apart from this, they have no effects in proving our main theorems (i.e., Theorem 1 and 2). These two assumptions hold by simply adding strongly convex regularizers to $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$ or $\tilde{f}_t(\mathbf{W})$. For example, we can add Tikhonov regularizer $(\nu_1 \|\mathbf{h}\|_2^2 + \nu_2 \|\mathbf{r}\|_2^2)/2$ (for some positive ν_1, ν_2) to $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$, then in such case $m_1 = \min(\nu_1, \nu_2)$. Adding such regularizers can be regarded as a way to promote smoothness and avoid over-fitting on \mathbf{W} , \mathbf{h} and \mathbf{r} (see Section VI). Moreover, including the regularizers in $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$ or $\tilde{f}_t(\mathbf{W})$ will not alter any steps in our analyses. In terms of our algorithms, for small regularization weights (e.g., $\nu_1, \nu_2 \ll 1$), such regularizers will only slightly alter steps (36) and (37) in Algorithm 1. For example, in (36), $\tilde{\ell}(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}, \mathbf{r})$ will become

$$\begin{aligned} \tilde{\ell}'(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}, \mathbf{r}) \triangleq & \frac{1}{2} \|\mathbf{v} - \mathbf{W}\mathbf{h} - \mathbf{r}\|_2^2 + \frac{\nu_1}{2} \|\mathbf{h}\|_2^2 \\ & + \lambda \|\mathbf{r}\|_1 + \frac{\nu_2}{2} \|\mathbf{r}\|_2^2. \end{aligned} \quad (39)$$

Thus for simplicity we omit such regularizers in the algorithm. However, both of our solvers (OPGD and OADMM) can be straightforwardly adapted to the regularized case.

Finally, we define the notion of *stationary point* of a constrained optimization problem.

Definition 1. Given a real Hilbert space \mathcal{X} , a differentiable function $g : \mathcal{X} \rightarrow \mathbb{R}$ and a set $\mathcal{K} \subseteq \mathcal{X}$, $x_0 \in \mathcal{K}$ is a stationary point of $\min_{x \in \mathcal{K}} g(x)$ if $\langle \nabla g(x_0), x - x_0 \rangle \geq 0$, for all $x \in \mathcal{K}$.

B. Main Results and Key Lemmas

In this section we present our main results (Theorem 1 and 2) and key lemmas to prove these results (Lemma 1 and 2). We only show proof sketches here. Detailed proofs are deferred to Section S-3 to S-8 in the supplemental material.

Theorem 1 (Almost sure convergence of the sequence of objective values). *In Algorithm 1, the nonnegative stochastic process $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$ converges a.s.. Furthermore, $\{f(\mathbf{W}_t)\}_{t \geq 1}$ converges to the same almost sure limit as $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$.*

Proof Sketch. The proof proceeds in two major steps. First, making use of the quasi-martingale convergence theorem [41, Theorem 9.4 & Proposition 9.5] (see Lemma S-8), we prove $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$ converges a.s. by showing the series of the expected positive variations of this process converges. A key step in proving the convergence of this series is to bound each summand of the series by Donsker's theorem (see Lemma S-9). To invoke this theorem, we show the class of measurable functions $\{\ell(\cdot, \mathbf{W}) : \mathbf{W} \in \mathcal{C}\}$ is \mathbb{P} -Donsker [40] using the Lipschitz continuity of $\ell(\mathbf{v}, \cdot)$ on \mathcal{C} (see Lemma 1) and [40, Example 19.7] (see Lemma S-11). Second, we show

$$f_t(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \xrightarrow{\text{a.s.}} 0. \quad (40)$$

by showing $\sum_{t=1}^{\infty} \frac{\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1}$ converges a.s.. Define $b_t \triangleq \tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t)$. Using the Lipschitz continuity of both \tilde{f}_t and f_t on \mathcal{C} and Lemma 2, we can show $|b_{t+1} - b_t| = O(1/t)$ a.s.. Now we invoke [44, Lemma 8] (see Lemma S-12) to conclude $f_t(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \xrightarrow{\text{a.s.}} 0$. By Lemma 1 and [40, Example 19.7] (see Lemma S-11), we can show the class of measurable functions $\{\ell(\cdot, \mathbf{W}) : \mathbf{W} \in \mathcal{C}\}$ is \mathbb{P} -Glivenko-Cantelli [40]. By the Glivenko-Cantelli theorem (see Lemma S-10) we have

$$\sup_{\mathbf{W} \in \mathcal{C}} |f_t(\mathbf{W}) - f(\mathbf{W})| \xrightarrow{\text{a.s.}} 0. \quad (41)$$

In particular,

$$f_t(\mathbf{W}_t) - f(\mathbf{W}_t) \xrightarrow{\text{a.s.}} 0. \quad (42)$$

Finally, (40) and (42) together imply that $\{f(\mathbf{W}_t)\}_{t \geq 1}$ converges to the same almost sure limit as $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$. \square

Theorem 2 (Almost sure convergence of the sequence of dictionaries). *The stochastic process $\{\mathbf{W}_t\}_{t \geq 1}$ converges to the set of stationary points of (8) a.s..⁴*

Proof Sketch. Fix a realization $\{\mathbf{v}_t\}_{t \geq 1}$ such that $\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t) \rightarrow 0$,⁵ and generate $\{\mathbf{W}_t\}_{t \geq 1}$ according to

⁴Given a metric space (\mathcal{X}, d) , a sequence (x_n) in \mathcal{X} is said to converge to a set $\mathcal{A} \subseteq \mathcal{X}$ if $\lim_{n \rightarrow \infty} \inf_{a \in \mathcal{A}} d(x_n, a) = 0$.

⁵The set of all such realizations have probability one by (40).

Algorithm 1. Then it suffices to show every subsequential limit of $\{\mathbf{W}_t\}_{t \geq 1}$ is a stationary point of (8). The compactness of \mathcal{C} enables us to find a convergent subsequence $\{\mathbf{W}_{t_m}\}_{m \geq 1}$. Furthermore, it is possible to find a further subsequence of $\{t_m\}_{m \geq 1}$, $\{k \geq 1\}$ such that all $\{\mathbf{A}_{t_k}\}_{k \geq 1}$, $\{\mathbf{B}_{t_k}\}_{k \geq 1}$ and $\{f_{t_k}(\mathbf{0})\}_{k \geq 1}$ converge. We focus on $\{\mathbf{W}_{t_k}\}_{k \geq 1}$ hereafter and denote its limit by $\overline{\mathbf{W}}$. Our goal is to show for any $\mathbf{W} \in \mathcal{C}$, the directional derivative $\langle \nabla f(\overline{\mathbf{W}}), \mathbf{W} - \overline{\mathbf{W}} \rangle \geq 0$. First, we show the sequence of differentiable functions $\{\tilde{f}_{t_k}\}_{k \geq 1}$ converges uniformly to a differentiable function \tilde{f} . By (41) we also have that $\{f_{t_k}\}_{k \geq 1}$ converges uniformly to f . Denote $g_t \triangleq \tilde{f}_t - f_t$, we have $\{g_{t_k}\}_{k \geq 1}$ converges uniformly to $\tilde{g} = \tilde{f} - f$. By Lemma 1, \tilde{f} is differentiable so \tilde{g} is differentiable on \mathcal{C} . Next, we show for any $\mathbf{W} \in \mathcal{C}$, $\langle \nabla \tilde{f}(\overline{\mathbf{W}}), \mathbf{W} - \overline{\mathbf{W}} \rangle \geq 0$ by showing $\overline{\mathbf{W}}$ is a global minimizer of \tilde{f} . Then we show $\nabla \tilde{g}(\overline{\mathbf{W}}) = \mathbf{0}$ using the first-order Taylor expansion. These results suffice to imply $\langle \nabla f(\overline{\mathbf{W}}), \mathbf{W} - \overline{\mathbf{W}} \rangle \geq 0$ for any $\mathbf{W} \in \mathcal{C}$. \square

Remark 5. Due to the nonconvexity of the canonical NMF problem, finding global optima of (2) is in general out-of-reach. Indeed, [68] shows (2) is NP-hard. Therefore in the literature [58], [69], convergence to the set of stationary points (see Definition 1) of (2) has been studied instead. In the online setting, it is even harder to show that the sequence of dictionaries $\{\mathbf{W}_t\}_{t \geq 1}$ converges to the global optima of (8) (in some probabilistic sense). Thus Theorem 2 only states the convergence result with respect to the stationary point of (8), which has been the state-of-the-art in the literature of online matrix factorization [37], [44]. On the other hand, we notice that under certain assumptions on the data matrix \mathbf{V} , e.g., the separability condition proposed in [70], exact NMF algorithms have been proposed in [71]–[73]. The optimization techniques therein are discrete (and combinatorial) in nature, and vastly differ from the continuous optimization techniques typically employed in solving the canonical NMF problem. In addition, some heuristics for obtaining exact NMF using continuous optimization techniques have been proposed in [74]. Nevertheless, developing exact NMF algorithms in the online setting could be an interesting future research direction.

The following two lemmas are used in the proof of Theorem 1 and 2. In particular, Lemma 1 states that the loss functions ℓ and f respectively defined in (5) and (8) satisfy some regularity conditions. Lemma 2 then bounds the variations in the stochastic process $\{\mathbf{W}_t\}_{t \geq 1}$.

Lemma 1 (Regularity properties of ℓ and f). *The loss functions ℓ and f satisfy the following properties*

- 1) $\ell(\mathbf{v}, \mathbf{W})$ is differentiable on $\mathcal{V} \times \mathcal{C}$ and $\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$ is continuous on $\mathcal{V} \times \mathcal{C}$. Moreover, for all $\mathbf{v} \in \mathcal{V}$, $\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$ is Lipschitz on \mathcal{C} with Lipschitz constant independent of \mathbf{v} .
- 2) The expected loss function $f(\mathbf{W})$ is differentiable on \mathcal{C} and $\nabla f(\mathbf{W})$ is Lipschitz on \mathcal{C} .

Proof Sketch. The regularity properties of ℓ on $\mathcal{V} \times \mathcal{C}$ originate from the regularity properties of $\tilde{\ell}$ on $\mathcal{V} \times \mathcal{C}$. Since Assumption 2 ensures the minimizer for (38),

$(\mathbf{h}^*(\mathbf{v}, \mathbf{W}), \mathbf{r}^*(\mathbf{v}, \mathbf{W}))$ is unique for any $(\mathbf{v}, \mathbf{W}) \in \mathcal{V} \times \mathcal{C}$, we can invoke Danskin’s theorem (see Lemma S-5) to guarantee the differentiability of ℓ . Moreover, we invoke the maximum theorem (see Lemma S-6) to ensure the continuity of $(\mathbf{h}^*(\mathbf{v}, \mathbf{W}), \mathbf{r}^*(\mathbf{v}, \mathbf{W}))$ on $\mathcal{V} \times \mathcal{C}$. These results together imply that ℓ satisfies the desired regularity properties. The regularity properties of f on \mathcal{C} hinges upon the regularity properties of ℓ . Indeed, by Leibniz integral rule (see Lemma S-7), we can show f is differentiable and

$$\nabla f(\mathbf{W}) = \mathbb{E}_{\mathbf{v}}[\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})]. \quad (43)$$

Hence the Lipschitz continuity of ∇f follows naturally from the Lipschitz continuity of $\mathbf{W} \mapsto \nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$. \square

Lemma 2. *In Algorithm 1, $\|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F = O(1/t)$ a.s..*

Proof Sketch. The upper bound for $\|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F$ results from the strong convexity and Lipschitz continuity of \tilde{f}_t . Specifically, they together imply an order difference between the upper and lower bounds of $\tilde{f}_t(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t)$, i.e.,

$$\frac{m_2}{2} \|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F^2 \leq \tilde{f}_t(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t) \quad (44)$$

$$\leq c_t \|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F, \quad (45)$$

where c_t is only dependent on t . Some calculations reveal that $c_t = O(1/t)$ a.s.. \square

C. Discussions

We highlight the distinctions between our convergence analyses and those in the previous works. Our first main result (Theorem 1) is somewhat analogous to the results in [16], [36], [37], [44]. However, the stochastic optimization of (8) in [16] is based on robust stochastic approximations [75], a very different approach from ours (see Section IV). This leads to significant differences in the analyses. For example, at time t , their dictionary update step does not necessarily minimize $\tilde{f}_t(\mathbf{W})$, a fact that we heavily leverage. The rest of the works fall under a similar framework as ours. However, in [44], a general sparse coding problem is considered. Thus, they assume a sufficient condition ensuring a unique solution in the Lasso-like sparse coding step holds. This assumption adds certain complications in proving ∇f is Lipschitz. Because our problem setting is different, we avoid these issues. Thus we are able to provide a more succinct proof, despite having some additional features, e.g., the presence of the outlier vector \mathbf{r} . The most closely related works are [36] and [37], both of which consider the online robust PCA problems but with different loss functions. However, the nonnegativity constraints on \mathbf{W} and \mathbf{h} and box constraints on \mathbf{r} distinguish our proof from theirs. For our second main result (Theorem 2), we note that in previous works [36], [37], seemingly stronger results have been shown. Specifically, these works manage to show $\{\mathbf{W}_t\}_{t \geq 1}$ converges, either to *one* stationary point [37] or to the global minimizer [36] of f . However, due to different

problem settings, their proof techniques cannot be easily adapted to our problem setting. Inspired by [60], we instead manage to characterize the subsequential limits of $\{\mathbf{W}_t\}_{t \geq 1}$ generated by almost sure realizations of $\{\mathbf{v}_t\}_{t \geq 1}$. This result is almost as good as the results in the abovementioned previous works since it implies that there exists a stationary point of f such that its neighborhood contains infinitely many points of the sequence almost surely.

VI. EXTENSIONS

In Section III, we considered a basic formulation of the online NMF problem with outliers. Here, we extend this formulation to a wider class of problem settings, with different constraints and regularizers on \mathbf{W} , \mathbf{h} and \mathbf{r} . We also discuss how to adapt our algorithms and analyses to these settings.

A. Different Constraints on \mathbf{W} and \mathbf{r}

If the constraint set \mathcal{C} changes, in terms of both solvers in Section IV, only (21) and (34), the updates involving projections onto \mathcal{C} , need to be changed. In terms of the analyses, using a similar argument as in Section V-A, we can show the (unique) optimal solution $\mathbf{W}^*(\mathbf{A}_t, \mathbf{B}_t)$ in (12) is bounded⁶ regardless of boundedness of \mathcal{C} . Thus, the optimization problem (12) can always be restricted to a compact set. The rest of the analyses proceeds as usual. Therefore, in the following, we only discuss the (convex) variants of \mathcal{C} and the associated (Euclidean) projection methods.

1) The nonnegative orthant $\mathcal{C}_1 \triangleq \mathbb{R}_+^{F \times K}$. This constraint is the most basic yet the most widely used constraint in NMF, although it does not well control the scale of \mathbf{W} . The projection onto the nonnegative orthant involves simple entrywise thresholding.

2) The probability simplex $\mathcal{C}_2 \triangleq \{\mathbf{W} \in \mathbb{R}_+^{F \times K} \mid \|\mathbf{W}_{:i}\|_1 = 1, \forall i \in [K]\}$ [16]. Different from \mathcal{C} , \mathcal{C}_2 also prevents some columns of \mathbf{W} from having arbitrarily small norms. Efficient projection algorithms onto the probability simplex have been well studied. See [76], [77] for details.

3) The nonnegative elastic-net ball [78]–[80]: $\mathcal{C}_3 \triangleq \{\mathbf{W} \in \mathbb{R}_+^{F \times K} \mid \gamma_1 \|\mathbf{W}_{:i}\|_1 + \gamma_2/2 \|\mathbf{W}_{:i}\|_2^2 \leq 1, \forall i \in [K]\}$, where both γ_1 and γ_2 are *nonnegative*. \mathcal{C}_3 is a general constraint set since it subsumes both \mathcal{C} (the nonnegative ℓ_2 ball) and the nonnegative ℓ_1 ball. Compared to \mathcal{C} , this constraint encourages sparsity on the basis matrix \mathbf{W} , thus leading to better interpretability on the basis vectors. Efficient algorithms for projection onto the elastic-net ball have been proposed in [44], [81].

For the outlier vector \mathbf{r} , the nonnegativity constraint can be added to \mathcal{R} to model (bounded) nonnegative outliers, so the new constraint $\mathcal{R}' \triangleq \{\mathbf{r} \in \mathbb{R}_+^F \mid \|\mathbf{r}\|_\infty \leq M\}$. In such case $\|\mathbf{r}\|_1 = \mathbf{1}^T \mathbf{r}$ so (14) amounts to a quadratic minimization program with box constraints. Both the algorithms and the analyses can be easily adapted to such a simpler case.

⁶From (12), we observe that \mathbf{W}_t is a function of only \mathbf{A}_t and \mathbf{B}_t if Assumption 3 holds.

B. Regularizers for \mathbf{W} , \mathbf{h} and \mathbf{r}

We use $\lambda_i \geq 0$ to denote the penalty parameters. For \mathbf{W} , the elastic-net regularization can be employed, i.e., $\lambda_1 \|\mathbf{W}\|_{1,1} + \lambda_2/2 \|\mathbf{W}\|_F^2$. This includes ℓ_1 and ℓ_2 regularizers on \mathbf{W} as special cases. As pointed out in [19], [82] and [78], ℓ_1 and ℓ_2 regularizers promote sparsity and smoothness on \mathbf{W} respectively. Because $\mathbf{W} \geq 0$, (12) with the elastic-net regularization is still quadratic. In terms of the analyses, Assumption 3 can be removed and the rest remains the same.

For \mathbf{h} , several regularizers can be used:

1) The Lasso regularizer $\lambda_3 \|\mathbf{h}\|_1$. It induces sparsity on \mathbf{h} , hence the optimization problem (13) with such regularizer is termed *nonnegative sparse coding* [38], [80].

2) The Tikhonov regularizer: $\lambda_4/2 \|\mathbf{h}\|_2^2$. This regularizer induces smoothness and avoids over-fitting on \mathbf{h} . Both the ℓ_1 and ℓ_2 regularizers preserve the quadratic nature of (13).

3) The group Lasso regularizer $\lambda_5 \sum_\alpha \sqrt{\xi_\alpha} \|\mathbf{h}_\alpha\|_2$, where \mathbf{h}_α 's are (non-overlapping) subvectors of \mathbf{h} with lengths ξ_α . This regularizer induces sparsity among groups of entries in \mathbf{h} . Efficient algorithms to solve (13) under this regularization have been proposed in [82], [83].

4) The sparse group Lasso regularizer $\lambda_6(\nu \sum_\alpha \sqrt{\xi_\alpha} \|\mathbf{h}_\alpha\|_2 + (1-\nu) \|\mathbf{h}\|_1)$, where $\nu \in [0, 1]$. Compared with the group lasso regularizer, this one also encourages sparsity within each group. Efficient algorithms for solving (13) with this regularization have been discussed in [84], [85].

Since (13) with each regularizer above is still a convex program, the analyses remain unchanged. Similar regularizers on \mathbf{h} can be applied to \mathbf{r} , so for simplicity we omit the discussions here. Compared to \mathbf{h} , the only differences are that \mathbf{r} may be negative and bounded. However, standard algorithms in such case are available in the literature [84]–[86].

Remark 6. Some remarks are in order. First, for certain pairs of constraints and regularizers above (e.g., the elastic-net ball constraint of \mathbf{W} and the elastic-net regularization on \mathbf{W}), the optimization problems are equivalent (admit the same optimal solutions) for specific value pairs of (γ_i, λ_i) . However, the algorithms and analyses for these equivalent problems can be much different. Second, all the alternative constraint sets and regularizers discussed above are convex, since our analyses heavily leverage the (strong) convexity of $\tilde{\ell}$ and \tilde{f} (see Section V-A). It would be a meaningful future research direction to consider nonconvex constraints or regularizers.

VII. EXPERIMENTS

We present results of numerical experiments to demonstrate the computational efficiency and efficacy of our online algorithms (OPGD and OADMM). We first state the experimental setup. Next we introduce some heuristics used in our experiments and the choices of parameters. We show the efficiency of our algorithms by examining the convergence speeds of the objective functions. The efficacy of our algorithms is shown via the (parts-based) basis representations and three meaningful real-world applications—image denoising, shadow removal and foreground-background separation. All the experiments were run in 64-bit Matlab[®] (R2015b) on a machine with Intel[®] Core i7-4790 3.6 GHz CPU and 8 GB RAM.

A. Experimental Setup

The datasets used in the experiments include one synthetic dataset, two face datasets (including the CBCL dataset [4] and the YaleB dataset [87]) and two video sequences in the i2r dataset [88]. For each task, we compared the performances of our online algorithms against those of four other matrix factorization algorithms, including BPGD, BADMM, online robust PCA (ORPCA) [36] and online NMF (ONMF) [16]. Such comparisons allow us to show multiple advantages of the proposed algorithms. Specifically, compared with their batch counterparts (BPGD and BADMM), our online algorithms have much faster convergence speeds in terms of the objective values.⁷ Moreover, compared with ORPCA, we are able to learn parts-based basis representations. Finally, compared with ONMF, we can recover the underlying low-dimensional data structure with minimal effects of outliers.

B. Strategies in Practical Implementations

1) *Initializations*: For both online algorithms (OPGD and OADMM), the entries of initial dictionary \mathbf{W}_0 were randomly generated independently and identically from the standard uniform distribution $\mathcal{U}[0, 1]$. This distribution can certainly be of other forms, e.g., the half-normal distribution⁸ $\mathcal{HN}(0, 1)$. However, we observed in all the experiments that different initialization schemes of \mathbf{W}_0 led to similar results. At each time instant, similar initialization methods on \mathbf{h} and \mathbf{r} were also used in solving (36). While for solving (37), we use \mathbf{W}_{t-1} as the initialization to exploit the similarities between dictionaries in adjacent iterations. This approach led to improved computational efficiency in practice. For the batch algorithms (BPGD and BADMM), the initial dictionary and coefficient matrix were similarly initialized as \mathbf{W}_0 .

2) *Stopping Criteria*: For the optimization problem (36), we employed the same stopping criterion for both algorithms OPGD and OADMM. Specifically, at time t , we stopped (36) at iteration k if

$$\frac{\left| \tilde{\ell}(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}^k, \mathbf{r}^k) - \tilde{\ell}(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}^{k-1}, \mathbf{r}^{k-1}) \right|}{\tilde{\ell}(\mathbf{v}_t, \mathbf{W}_{t-1}, \mathbf{h}^{k-1}, \mathbf{r}^{k-1})} < \epsilon_{\text{th}}$$

or $k > \varkappa_{\text{max}}$, where $\epsilon_{\text{th}} > 0$ denotes the threshold for the relative decrease of objective values and $\varkappa_{\text{max}} \in \mathbb{N}$ denotes the maximum number of iterations. We set $\epsilon_{\text{th}} = 1 \times 10^{-3}$ and $\varkappa_{\text{max}} = 50$ for all the experiments. Similar stopping criterion applies to both OPGD and OADMM in solving (37), except in this case, we set $\epsilon_{\text{th}} = 1 \times 10^{-4}$ and $\varkappa_{\text{max}} = 200$ for better precision in updating the dictionary.

⁷ The problem formulations in the existing batch robust NMF algorithms [13], [26]–[32], [34], [35] are different from ours here. Specifically, most of the algorithms do not use the $\ell_{1,1}$ regularization to remove the effect of outliers. Moreover, the constraints on \mathbf{W} and \mathbf{R} are simpler in their formulations. All these differences contribute to different results in the experiments. Thus for the purpose of fair comparisons, we chose to compare our online algorithms with their batch counterparts.

⁸ $\mathcal{HN}(\sigma^2)$ denotes the half-normal distribution with scale parameter σ^2 , i.e., $\mathcal{HN}(y; \sigma^2) = \frac{\sqrt{2}}{\sigma\sqrt{\pi}} \exp(-\frac{y^2}{2\sigma^2})$ for $y \geq 0$. If $Y \sim \mathcal{N}(0, \sigma^2)$, $|Y|$ is half-normal with scale parameter σ^2 .

3) *Heuristics*: We now describe some heuristics employed in the implementations of our online algorithms. We remark that all the heuristics used in this work are common in the literature of online matrix factorization [15], [44]. First, since the sample size of the *benchmarking* datasets in real world may not be large enough, we replicate it p times to aggregate its size. Each (aggregated) dataset is represented as an $F \times N$ matrix, where F equals the ambient dimension of the data⁹ and N equals p times the size of the original dataset. Next, for the aggregated image dataset, we randomly shuffle the images in it to simulate an i.i.d. data stream. We also normalize each data sample so that it has unit maximum pixel value. The last heuristic is called the mini-batch extension, namely at a time instant, we process τ images altogether ($\tau \geq 1$). This amounts to sampling τ i.i.d. samples from \mathbb{P} at each time.

C. Parameter Settings

We discuss how to set the values of the important parameters in our online algorithms. These parameters include the regularization parameter λ in (5), the constraint parameter M in the definition of \mathcal{R} , the mini-batch size τ , the latent data dimensionality K , the penalty parameters $\{\rho_i\}_{i \in [3]}$ in OADMM, and the step sizes $\bar{\kappa}$ and $\{\tilde{\kappa}_t\}_{t \geq 1}$ in OPGD. First, we set $\lambda = 1/\sqrt{F}$, following the convention in the literature [36], [37], [44]. Since each image has unit maximum pixel value, it is straightforward to set $M = 1$. For the mini-batch size τ , we propose a rule-of-thumb, that is to choose $\tau \leq 4 \times 10^{-4}N$, e.g., $\tau = 5 \times 10^{-5}N$. The choice of this parameter involves a trade-off between the stability of the basis matrix and the convergence speed of the online algorithms. In general this parameter is data-dependent so there is no principled way of choosing it in the literature [16], [44]. Therefore, our approach presents a heuristic way to partially resolve this issue. Regarding the latent data dimensionality K , there are some works [89], [90] that describe how to choose this parameter from a Bayesian perspective. However, they involve complex Bayesian modeling and intensive computations. Thus for efficiency we set $K = 49$. Finally we discuss the choice of the parameters specific to each online algorithm. We first reduce the number of parameters in both algorithms by setting $\rho_i = \rho$, for any $i \in [3]$ in OADMM¹⁰ and $\bar{\kappa} = \tilde{\kappa}_t = \kappa$, for any $t \geq 1$ in OPGD. Then we set $\rho = 1$ and $\kappa = 0.7$. The above parameter setting will be referred as *the canonical parameter setting* in the sequel. Later we will show the convergence speeds of both our online algorithms are insensitive to the parameters that are set to fixed values (including τ , K , ρ and κ) within wide ranges on both the synthetic and CBCL datasets.

D. Convergence Speeds

The convergence speeds of the objective values of OPGD and OADMM are compared to the other four benchmarking algorithms (BPGD, BADMM, ORPCA and ONMF) on both the synthetic and CBCL datasets. For better illustration,

⁹Each image (or video frame) is vectorized and stacked as a column of the data matrix, so F equals the number of pixels in the image (or video frame).

¹⁰See the literature on ADMM (e.g., [91]) for more sophisticated ways of choosing this parameter.

the comparison is divided into two parts. The first involves comparison between our online algorithms and their batch counterparts. The second involves comparison between our online algorithms to other online algorithms. For the first part, we use the surrogate loss function \hat{f}_t as a measure of convergence because i) $\hat{f}_t(\mathbf{W}_t)$ has the same a.s. limit as $f(\mathbf{W}_t)$ (as $t \rightarrow \infty$) and ii) with storage of the past statistics $\{\mathbf{v}_i, \mathbf{h}_i, \mathbf{r}_i\}_{i \in [t]}$, \hat{f}_t is easier to compute than the expected loss f or the empirical loss f_t . For the second part, since different online algorithms have different objective functions, we propose a heuristic and unified measure of convergence $\hat{f}_t: \mathcal{C} \rightarrow \mathbb{R}_+$, defined as

$$\hat{f}_t(\mathbf{W}) \triangleq \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{v}_i^o - \mathbf{W}\mathbf{h}_i\|_2^2, \quad (46)$$

where \mathbf{v}_i^o denotes the clean data sample (without outliers and observation noise) at time i . Loosely speaking, $\hat{f}_t(\mathbf{W}_t)$ can be interpreted as the averaged regret of reconstructing the clean data up to time t . In the following experiments, unless compared with other online algorithms, we always use \hat{f}_t as the objective function for our online algorithms.

1) *Generation of Synthetic Dataset*: The procedure to generate the synthetic dataset is described as follows. First, we generated $\mathbf{W}^o \in \mathbb{R}_+^{F \times K^o}$ such that $(w^o)_{ij} \stackrel{iid}{\sim} \mathcal{HN}(1/\sqrt{K^o})$, for any $(i, j) \in [F] \times [K^o]$, where K^o denotes the ground-truth latent dimension. Similarly, we generated $\mathbf{H}^o \in \mathbb{R}_+^{K^o \times N}$ such that $(h^o)_{ij} \stackrel{iid}{\sim} \mathcal{HN}(1/\sqrt{K^o})$, for any $(i, j) \in [K^o] \times [N]$. The (clean) data matrix $\mathbf{V}^o = \mathcal{P}_{\tilde{\mathcal{V}}}(\mathbf{W}^o\mathbf{H}^o)$, where $\tilde{\mathcal{V}} = [0, 1]^{F \times N}$. Note that the normalization (projection) $\mathcal{P}_{\tilde{\mathcal{V}}}$ preserves the boundedness of the data. Then we generated the outlier vectors $\{\mathbf{r}_i\}_{i \in [N]}$. First, we uniformly randomly selected a subset of $[N]$ with size $\lfloor \nu N \rfloor$ and denoted it as \mathcal{I} , where $\nu \in (0, 1)$ denotes the fraction of columns in \mathbf{V}^o to be contaminated by the outliers. For each $i \in \mathcal{I}$, we generated the outlier vector $\mathbf{r}_i \in [-1, 1]^F$ by first uniformly selecting a support set with cardinality $\lfloor \tilde{\nu} F \rfloor$, where $\tilde{\nu} \in (0, 1)$ denotes the outlier density in each data sample. Then the nonzero entries in \mathbf{r}_i were i.i.d. generated from the uniform distribution $\mathcal{U}[-1, 1]$. For each $i \notin \mathcal{I}$, we set $\mathbf{r}_i = \mathbf{0}$. Define the outlier matrix \mathbf{R}^o such that $(\mathbf{R}^o)_{:i} = \mathbf{r}_i$, for any $i \in [N]$. Then the (contaminated) data matrix $\mathbf{V} = \mathcal{P}_{\tilde{\mathcal{V}}}(\mathbf{V}^o + \mathbf{R}^o + \mathbf{N})$, where $\mathbf{N} \in \mathbb{R}^{F \times N}$ denotes the observation noise matrix with i.i.d. standard normal entries.

2) *Comparison to Other Online and Batch NMF Algorithms*: To make a fair comparison, all the algorithms were initialized with the same \mathbf{W}_0 for each parameter setting. Moreover, all online algorithms had the same mini-batch size τ . The *synthetic* data matrix \mathbf{V} was generated using the parameters values $F = 400$, $K^o = 49$, $N = 1 \times 10^5$, $\nu = 0.7$ and $\tilde{\nu} = 0.1$.

The convergence speeds of the aforementioned algorithms under the canonical parameter setting¹¹ are shown in Figure 1. From Figure 1(a), we observe our online algorithms converge much faster than their batch counterparts. From Figure 1(b), we observe the ORPCA algorithm converges slightly faster

than our online algorithms (the time difference is less than 1s). This is because ORPCA's formulation does not incorporate nonnegativity constraints (and magnitude constraints for the outliers) so the algorithm has fewer projection steps, leading to a lower computational complexity. However, we will show in Section VII-E that it fails to learn the parts-based representations due to the lack of nonnegativity constraints. Moreover, we also observe that the ONMF algorithm fails to converge because it is unable to handle the outliers. Thus the constant perturbation from the outliers in the data samples on the learned basis matrix keeps the algorithm from converging. Furthermore, from both subfigures, we observe the objective function of OADMM has a smoother decrease than that of OPGD. This is OADMM solves the optimization problems (9) and (12) in the dual space so it avoids the sharp transitions from plateaus to valleys in the primal space. Apart from this difference, *the convergence speeds of OADMM and OPGD are almost the same.*

3) *Insensitivity to Parameters*: We now examine the influences of various parameters, including τ , K , ρ and κ , on the convergence speeds of our online algorithms on the synthetic dataset. To do so, we vary one parameter at a time while keeping the other parameters fixed as in the canonical setting. We first vary τ from 5 to 40 in a log-scale, since our proposed heuristic rule indicates $\tau \leq 40$. Figure 2 shows that our rule-of-thumb works well since within its indicated range, different values of τ have very little effects on the convergence speeds of both our online algorithms. Next, we consider other values of the latent dimension K . In particular, we consider $K = 25$ and $K = 100$. Figure 3 shows that the convergence speeds of both our online and batch algorithms are insensitive to this parameter. Then we vary ρ (in both BADMM and OADMM) from 1 to 1000 on a log-scale. We can make two observations from Figure 4. First, as ρ increases, both ADMM-based algorithms converge slower. However, the change of the convergence speed of OADMM is much less drastic compared to BADMM. In fact, in less than 10s, all the OADMM algorithms (with different values of ρ) converge. This is because for OADMM, a large ρ only has a significant effect on its convergence speed in the initial transient phase. After the algorithm reaches its steady state (i.e., the basis matrix becomes stable), the effect of a large ρ becomes minimal in terms of solving both (9) and (12). Similar effects of the different values of κ on the convergence speeds of the PGD-based algorithms are observed in Figure 5. Together, Figure 4 and 5 reveal *another advantage* of our online algorithms. That is, our online algorithms have a significantly larger tolerance of suboptimal parameters than their batch counterparts.

4) *Effects of Proportion of Outliers*: We evaluate the performance of our online (and batch) algorithms (with the canonical parameter setting) on the synthetic dataset with a larger proportion of outliers. Specifically, in generating the synthetic data matrix \mathbf{V} , we keep F , K^o and N unchanged, but simultaneously increase ν and $\tilde{\nu}$. Figures 6(a) and 6(b) show the convergence results of our algorithms on the synthetic dataset with $(\nu, \tilde{\nu})$ equal to $(0.8, 0.2)$ and $(0.9, 0.3)$ respectively. From Figure 6, we observe that our online algorithms still converge fast and steadily (without fluctuations) to a

¹¹For BPGD and BADMM, the step sizes and penalty parameters had the same values as those of the online algorithms. For ORPCA and ONMF, we kept the parameter settings in the original works.

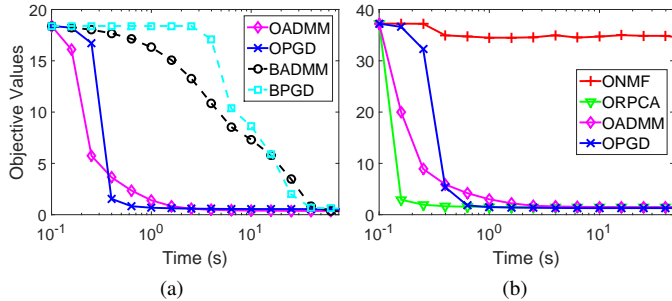


Fig. 1. The objective values (as a function of time) of (a) our online algorithms and their batch counterparts (b) our online algorithms and other online algorithms. The parameters are set according to the canonical setting.

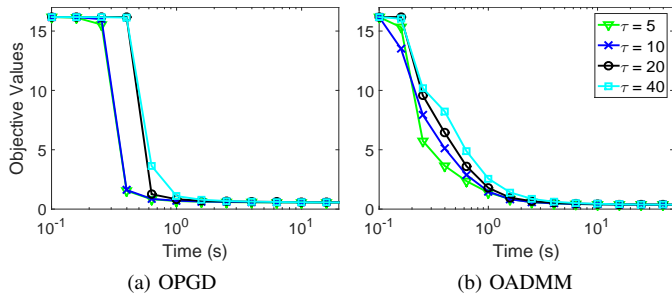


Fig. 2. The objective values of (a) OPGD and (b) OADMM for different values of τ . All the other parameters are set according to the canonical setting.

stationary point even though both ν and $\tilde{\nu}$ have increased.

5) *Experiments on the CBCL Dataset:* To generate the CBCL face dataset with outliers, we first replicated it by a factor $p = 50$ (so that the total number of data samples is of the order 10^5) and randomly permuted the aggregated dataset as described in Section VII-B. We then contaminated it with outliers in the same manner as for the synthetic dataset. However, we avoided adding observation noise since it had been already introduced by the image acquisition process.

We then conducted experiments on the contaminated face dataset with the same parameter settings as those on the synthetic data. In the interest of space, we defer the convergence results of all the online and batch algorithms to Section S-10. The results show that both our online algorithms demonstrate fast and steady convergences on this dataset. Moreover, the convergence speeds are insensitive to the key parameters (τ , K , ρ and κ). Therefore, in the subsequent experiments, we will only focus on the canonical parameter setting unless mentioned otherwise.

E. Basis Representations

We now examine the basis matrices learned by all the algorithms on the CBCL dataset with outliers as introduced in Section VII-D5. The parameters controlling the outlier density, ν and $\tilde{\nu}$ are set to 0.7 and 0.1 respectively. Figure 7 shows the basis representations learned by all the algorithms. From this figure, we observe that the basis images learned by ONMF have large residues of salt and pepper noise. Also, the basis images learned by ORPCA appear to be non-local and do not yield meaningful interpretations. In contrast, the basis images

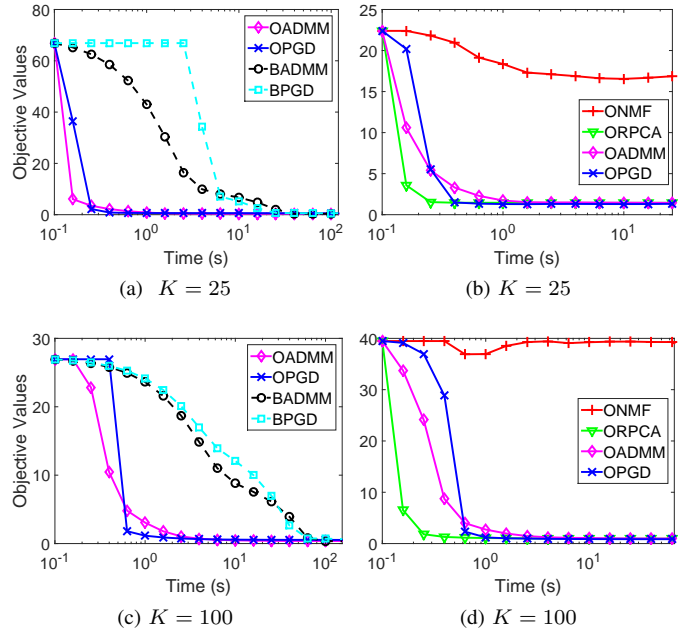


Fig. 3. The objective values (as a function of time) of all the algorithms for different values of K . In (a) and (b), $K = 25$. In (c) and (d), $K = 100$. All the other parameters are set according to the canonical setting.

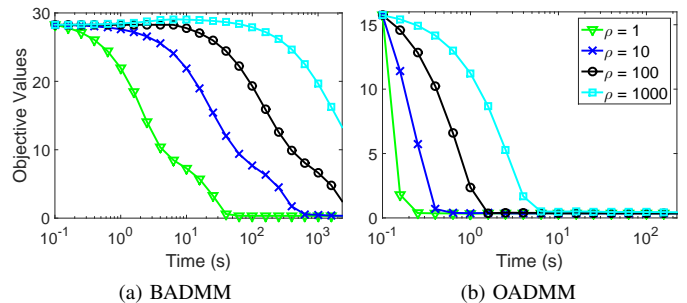


Fig. 4. The objective values of (a) BADMM and (b) OADMM for different values of ρ . All the other parameters are set according to the canonical setting.

learned by our online (and batch) algorithms are free of noise and much more local than those learned by ORPCA. We can easily observe facial parts from the learned basis. Also, the basis images learned by our four algorithms are of similar quality. To further enhance the sparsity of the basis images learned, one can add sparsity constraints to each column of the basis matrix, as done for example in [43]. However, we omit such constraints in our discussions. Another parameter affecting the sparsity of the learned basis images is the latent dimension K . In general, the larger K is, the sparser the basis images will be. Here we set $K = 49$, in accordance with the original setting in [3].

F. Application I: Image Denoising

A natural application that arises from the experiments on the CBCL dataset would be image denoising, in particular, removing the salt and pepper noise on images. The metric commonly used to measure the quality of the reconstructed images is peak signal-to-noise ratio (PSNR) [92]. A larger

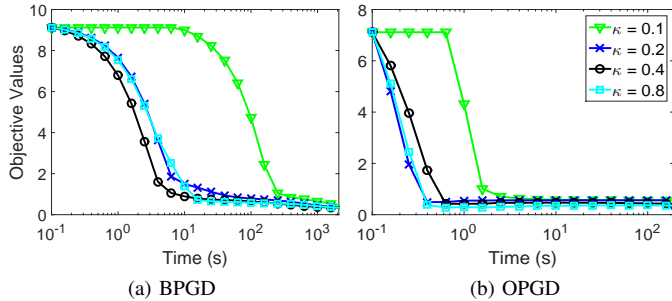


Fig. 5. The objective values of (a) BPGD and (b) OPGD for different values of κ . All the other parameters are set according to the canonical setting.

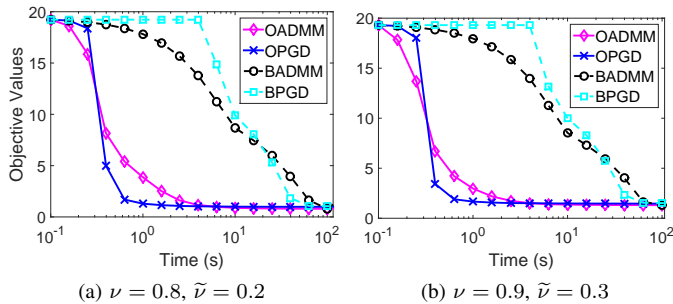


Fig. 6. The objective values of our online and batch algorithms on the synthetic dataset with a larger proportion of outliers.

value of PSNR indicates better quality of image denoising. With a slight abuse of definition, we apply this metric to all the recovered images. For the batch algorithms, we define

$$\text{PSNR} \triangleq -10 \log_{10} \left\{ \frac{\|\mathbf{V}^o - \widehat{\mathbf{W}}\widehat{\mathbf{H}}\|_F^2}{FN} \right\}, \quad (47)$$

where \mathbf{V}^o , $\widehat{\mathbf{W}}$ and $\widehat{\mathbf{H}}$ denote the matrix of the clean images, estimated basis matrix and estimated coefficient matrix respectively. For the online algorithms, we instead define PSNR in terms of \mathbf{W}_N (the final dictionary output by Algorithm 1) and the past statistics $\{\mathbf{h}_i\}_{i \in [N]}$

$$\text{PSNR} \triangleq -10 \log_{10} \left\{ \sum_{i=1}^N \|\mathbf{v}_i^o - \mathbf{W}_N \mathbf{h}_i\|_2^2 / (FN) \right\} \quad (48)$$

$$\stackrel{c}{=} -10 \log_{10} \widehat{f}_N(\mathbf{W}_N). \quad (49)$$

In other words, a low averaged regret will result in a high PSNR. Table I shows the image denoising results of all the algorithms on the CBCL face dataset. Here the settings 1, 2 and 3 represent different densities of the salt and pepper noise. Specifically, these settings correspond to $(\nu, \tilde{\nu})$ equal to $(0.7, 0.1)$, $(0.8, 0.2)$ and $(0.9, 0.3)$ respectively. All the results were obtained using ten random initializations of \mathbf{W}_0 . From Table I(a), we observe that for all the three settings, in terms of PSNRs, our online algorithms only slightly underperform their batch counterparts, but slightly outperform ORPCA and greatly outperform ONMF. From Table I(b), we observe our online algorithms only have slightly longer running times than ORPCA, but are significantly faster than the rest ones. Thus in terms of the trade-off between the computational efficiency and the quality of image denoising, our online algorithms

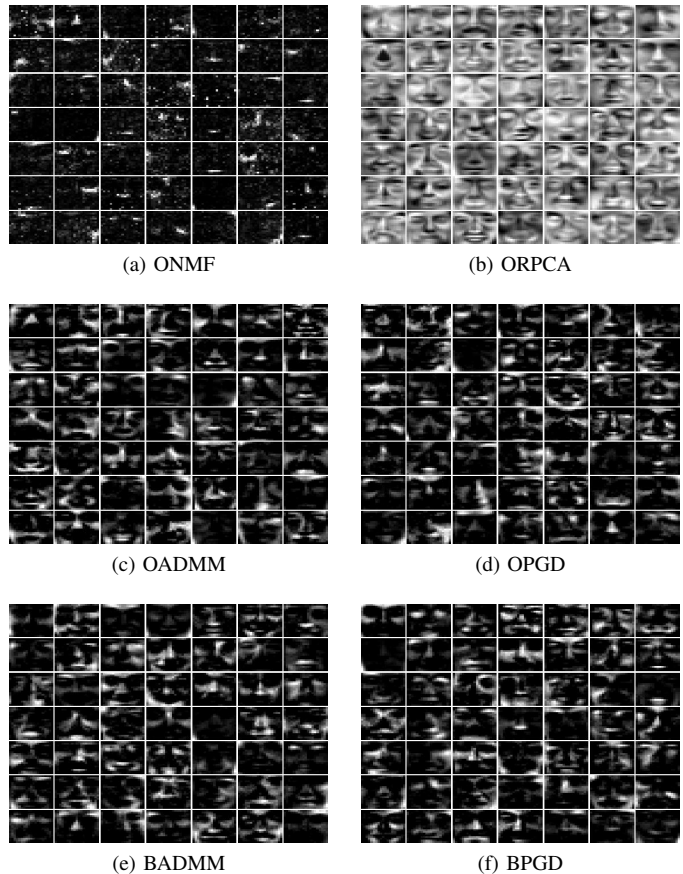


Fig. 7. Basis representations learned by all the algorithms ($\nu = 0.7$, $\tilde{\nu} = 0.1$). All the parameters are in the canonical setting ($K = 49$).

TABLE I
AVERAGE PSNRs (IN DB) AND RUNNING TIMES (IN SECONDS) WITH STANDARD DEVIATIONS OF ALL THE ALGORITHMS ON THE CBCL FACE DATASET WITH DIFFERENT NOISE DENSITY ($K = 49$).

	Setting 1	Setting 2	Setting 3
OADM	11.37 \pm 0.02	11.35 \pm 0.15	11.33 \pm 0.11
OPGD	11.48 \pm 0.10	11.47 \pm 0.05	11.39 \pm 0.03
BADM	11.53 \pm 0.19	11.51 \pm 0.10	11.48 \pm 0.03
BPGD	11.56 \pm 0.07	11.52 \pm 0.14	11.48 \pm 0.05
ONMF	5.99 \pm 0.12	5.97 \pm 0.18	5.95 \pm 0.17
ORPCA	11.25 \pm 0.04	11.24 \pm 0.19	11.23 \pm 0.05

(a) PSNRs

	Setting 1	Setting 2	Setting 3
OADM	416.68 \pm 3.96	422.35 \pm 3.28	425.83 \pm 4.25
OPGD	435.32 \pm 4.80	449.25 \pm 3.18	458.58 \pm 4.67
BADM	1000.45 \pm 10.18	1190.55 \pm 5.88	1245.39 \pm 10.59
BPGD	1134.89 \pm 11.37	1185.84 \pm 9.83	1275.48 \pm 9.48
ONMF	2385.93 \pm 11.15	2589.38 \pm 15.57	2695.47 \pm 14.15
ORPCA	368.35 \pm 3.23	389.59 \pm 3.49	399.85 \pm 4.12

(b) Running times

achieve comparable performances with ORPCA but are much better than the rest algorithms. Also, in terms of PSNRs and running times, no significant differences can be observed between our online algorithms. Similar results were observed when $K = 25$ and $K = 100$. See Section S-10 for details.

TABLE II
AVERAGE RUNNING TIMES (IN SECONDS) OF ALL THE ALGORITHMS ON
SUBJECT No.2.

Algorithms	Time (s)	Algorithms	Time (s)
ONMF	436.11 ± 10.85	ORPCA	16.38 ± 2.43
OADMM	51.12 ± 2.69	OPGD	62.34 ± 2.75
BADMM	175.52 ± 7.84	BPGD	212.71 ± 6.33

G. Application II: Shadow Removal

We evaluated the performances of all the online and batch algorithms on removing shadows in the face images in the `YaleB` dataset. It is well-known from the image processing literature that the shadows in an image can be regarded as inherent outliers. Therefore, the shadow removal task serves as another meaningful application of our online (and batch) algorithms. We first briefly describe the experiment procedure. For each subject in the `YaleB` dataset, we aggregated and randomly permuted his/her images. We set the aggregation factor p to 50 in consistency with the previous experiments. We then regarded the resulting data as the input to all the algorithms. Finally, we reconstructed the images in a similar way as in Section VII-F. The experiments were run using ten random initializations of \mathbf{W}_0 . Figure 8 shows some (randomly sampled) reconstructed images on subjects No.2 and No.8 (with one initialization of \mathbf{W}_0). From this figure, we observe that overall, our online algorithms perform almost as well as their batch counterparts, except for small artifacts (e.g., salt noise) in the recovered images by our online algorithms. We also observe that the other two online algorithms are inferior to our online algorithms. Specifically, ORPCA has more prominent artifacts (e.g., salt noise with larger density) in the recovered images. ONMF in most cases either fails to remove shadow or causes large distortions to the original images. The results by other initializations of \mathbf{W}_0 are similar to those shown in Figure 8. Table II shows the running times of all algorithms on the images of subject No.2. The average running times of these algorithms on other subjects are similar to those on subject No.2. This table, together with Figure 8, suggests that our online algorithms achieve the best trade-off between the computational efficiency and the quality of shadow removal. Also, both of our online algorithms have similar performances on the shadow removal task.

H. Application III: Foreground-Background Separation

Finally, we also evaluated the performance of all the online and batch algorithms on the task of foreground-background separation, which is important in video surveillance [88], [93]. Since the foreground objects in each video frame generally only occupy a small fraction of pixels, they have been considered as outliers in the literature [93], [94]. Therefore, our online (and batch) algorithms can be applied to estimate the background scene, as well as learn the foreground objects. In this section, we consider two video sequences, `Hall` and `Escalator` from the `i2r` dataset [88]. Each video sequence consists of 200 video frames and the resolutions of the frames in `Hall` and `Escalator` are 144×176 and 130×160

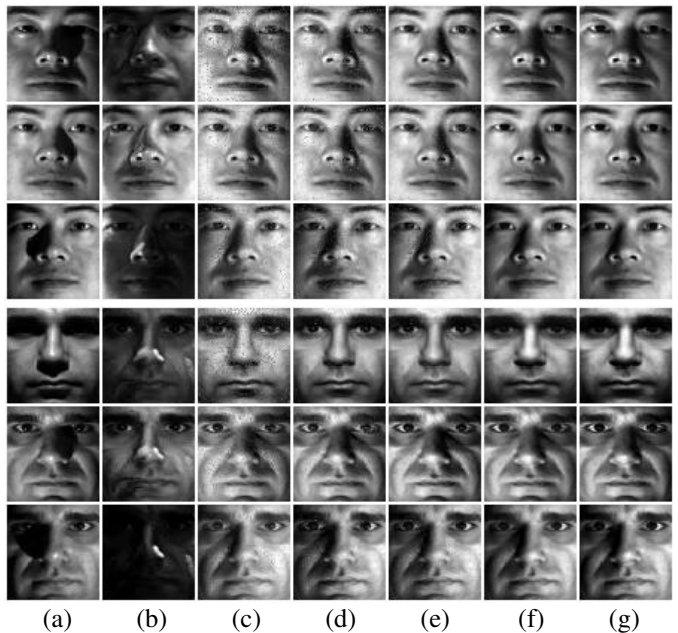


Fig. 8. Results of shadow removal by all the online and batch algorithms on subjects No.2 (upper) and No.8 (lower) in the `YaleB` face dataset. The labels (a) to (g) denote the original images and results produced by ONMF, ORPCA, OADMM, OPGD, BADMM, BPGD respectively.

respectively.¹² We set $p = 10$ for storage space considerations of the batch algorithms. We also repeated the experiments using ten random initializations of \mathbf{W}_0 .

The average running times over the ten initializations of \mathbf{W}_0 (with standard deviations) on the two video sequences are shown in Table III. We notice that consistent messages are delivered by Table III as those in Table I(b) and Table II. Namely, the running times of our online algorithms are slightly longer than those of ORPCA but greatly shorter than the rest algorithms. Figure 9 shows some (randomly sampled) background scenes and foreground objects separated by each algorithm. For each frame, the background was reconstructed using the methods introduced in Section VII-F, and the foreground was directly recovered from the corresponding column in the (estimated) outlier matrix $\hat{\mathbf{R}}$. Since $\hat{\mathbf{R}}$ is absent in the formulation of ONMF, it is estimated using the difference between the original video frames and the recovered background scenes. From Figure 9, we observe that on both video sequences, our online algorithms are able to separate the foreground objects from the background fairly successfully, with unnoticeable residues on both foreground and background. Compared with the other algorithms, the separation results of our online algorithms are comparable to their batch counterparts and slightly better than ORPCA on the `Hall` sequence (with less residues in the recovered foreground). Due to the lack of robustness, the background scenes recovered by ONMF appear to be dark and the foreground objects cannot be separated from the background. Thus again, in terms of the trade-off between the computational efficiency and the quality of foreground-background separation, our online algorithms achieve the best

¹²For simplicity we converted the color video frames to gray-scale using the built-in Matlab function `rgb2gray`.

TABLE III
AVERAGE RUNNING TIMES (IN SECONDS) OF ALL THE ALGORITHMS ON VIDEO SEQUENCES (A) Hall AND (B) Escalator.

Algorithms	Time (s)	Algorithms	Time (s)
ONMF	1525.57 ± 14.43	ORPCA	166.85 ± 6.09
OADMM	172.29 ± 5.64	OPGD	178.46 ± 7.57
BADMM	1280.06 ± 13.57	BPGD	1167.95 ± 13.38

(a)

Algorithms	Time (s)	Algorithms	Time (s)
ONMF	1324.86 ± 11.45	ORPCA	160.85 ± 5.03
OADMM	166.83 ± 5.64	OPGD	170.46 ± 5.91
BADMM	898.47 ± 8.57	BPGD	867.41 ± 9.35

(b)

performances. Also, both of our online algorithms have similar performances on the foreground-background separation task.

VIII. CONCLUSION AND FUTURE WORK

In this paper, we have developed online algorithms for NMF where the data samples are contaminated by outliers. We have shown that the proposed class of algorithms is robust to sparsely corrupted data samples and performs efficiently on large datasets. Finally, we have proved almost sure convergence guarantees of the objective function as well as the sequence of basis matrices generated by the algorithms.

We hope to pursue the following three directions for further research. First, it would be interesting to extend the convergence analyses given i.i.d. data sequences to weakly dependent data sequences (e.g., martingale or auto-regressive processes). This is because real data (e.g., video and audio recordings) have correlations among adjacent data samples. Second, it would be meaningful to consider the case where the fit-to-data term between \mathbf{v} and $\mathbf{W}\mathbf{h} + \mathbf{r}$, $d(\mathbf{v} \parallel \mathbf{W}\mathbf{h} + \mathbf{r})$ is not the squared ℓ_2 loss, e.g., the β -divergence [95], [96]. This is because in many scenarios, the observation noise is not Gaussian [95] so other types of loss functions need to be used. Finally, as mentioned in Remark 6, it would be meaningful to consider nonconvex constraint sets and/or regularizers on \mathbf{W} , \mathbf{h} or \mathbf{r} . This is because the nonconvex constraints and regularizers often appear in real applications [97].

Acknowledgment: The authors are grateful for the helpful clarifications and suggestions from Julien Mairal.

REFERENCES

- [1] R. Zhao and V. Y. F. Tan, "Online nonnegative matrix factorization with outliers," in *Proc. ICASSP*, Shanghai, China, Mar. 2016, pp. 2662–2666.
- [2] S. Tsuge, M. Shishibori, S. Kuroiwa, and K. Kita, "Dimensionality reduction using non-negative matrix factorization for information retrieval," in *Proc. SMC*, vol. 2, Tucson, Arizona, USA, Oct. 2001, pp. 960–965.
- [3] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, pp. 788–791, October 1999.
- [4] —, "Algorithms for non-negative matrix factorization," in *Proc. NIPS*, Denver, USA, Dec. 2000, pp. 556–562.
- [5] J. Kim and H. Park, "Toward faster nonnegative matrix factorization: A new algorithm and comparisons," in *Proc. ICDM*, Pisa, Italy, Dec. 2008, pp. 353–362.

- [6] C.-J. Lin, "Projected gradient methods for nonnegative matrix factorization," *Neural Comput.*, vol. 19, no. 10, pp. 2756–2779, Oct. 2007.
- [7] H. Kim and H. Park, "Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method," *SIAM J. Matrix Anal. A.*, vol. 30, no. 2, pp. 713–730, 2008.
- [8] Y. Xu, W. Yin, Z. Wen, and Y. Zhang, "An alternating direction algorithm for matrix completion with nonnegative factors," *Frontiers Math. China*, vol. 7, pp. 365–384, 2012.
- [9] C. Ding, X. He, and H. D. Simon, "On the equivalence of nonnegative matrix factorization and spectral clustering," in *Proc. SDM*, Newport Beach, California, USA, Apr. 2005, pp. 606–610.
- [10] Y. Yuan, Y. Feng, and X. Lu, "Projection-based nmf for hyperspectral unmixing," *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, vol. 8, no. 6, pp. 2632–2643, Jun. 2015.
- [11] J. L. Durrieu, B. David, and G. Richard, "A musically motivated mid-level representation for pitch estimation and musical audio source separation," *IEEE J. Sel. Top. Signal Process.*, vol. 5, no. 6, pp. 1180–1191, Oct. 2011.
- [12] B. Cao, D. Shen, J.-T. Sun, X. Wang, Q. Yang, and Z. Chen, "Detect and track latent factors with online nonnegative matrix factorization," in *Proc. IJCAI*, Hyderabad, India, Jan. 2007, pp. 2689–2694.
- [13] L. Zhang, Z. Chen, M. Zheng, and X. He, "Robust nonnegative matrix factorization," *Frontiers Elect. Electron. Eng. China*, vol. 6, no. 2, pp. 192–200, 2011.
- [14] Netflix, "Rad - outlier detection on big data," <http://techblog.netflix.com/2015/02/rad-outlier-detection-on-big-data.html>, 2015.
- [15] S. S. Bucak and B. Günsel, "Incremental subspace learning via nonnegative matrix factorization," *Pattern Recogn.*, vol. 42, no. 5, pp. 788–797, May 2009.
- [16] N. Guan, D. Tao, Z. Luo, and B. Yuan, "Online nonnegative matrix factorization with robust stochastic approximation," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 23, no. 7, pp. 1087–1099, Jul. 2012.
- [17] A. Lefèvre, F. Bach, and C. Févotte, "Online algorithms for nonnegative matrix factorization with the itakura-saito divergence," in *Proc. WASPAA*, New Paltz, New York, USA, Oct. 2011, pp. 313–316.
- [18] F. Wang, C. Tan, A. C. König, and P. Li, "Efficient document clustering via online nonnegative matrix factorizations," in *Proc. SDM*, Mesa, Arizona, USA, Apr. 2011, pp. 908–919.
- [19] Y. Wu, B. Shen, and H. Ling, "Visual tracking via online nonnegative matrix factorization," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 24, no. 3, pp. 374–383, Mar. 2014.
- [20] C. Liu, H. chih Yang, J. Fan, L.-W. He, and Y.-M. Wang, "Distributed nonnegative matrix factorization for web-scale dyadic data analysis on mapreduce," in *Proc. WWW*, Raleigh, North Carolina, USA, Apr. 2010, pp. 681–690.
- [21] R. Gemulla, P. J. Haas, Y. Sismanis, C. Teflioudi, and F. Makari, "Large-scale matrix factorization with distributed stochastic gradient descent," in *Proc. KDD*, San Diego, California, USA, Aug. 2011, pp. 69–77.
- [22] S. S. Du, Y. Liu, B. Chen, and L. Li, "Maxios: Large scale nonnegative matrix factorization for collaborative filtering," in *Proc. NIPS-DMLMC*, Montréal, Canada, Dec. 2014.
- [23] J. Chen, Z. J. Towfic, and A. H. Sayed, "Dictionary learning over distributed models," *IEEE Trans. Signal Process.*, vol. 63, no. 4, pp. 1001–1016, 2015.
- [24] N. Halko, P. G. Martinsson, and J. A. Tropp, "Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions," *SIAM Rev.*, vol. 53, no. 2, pp. 217–288, 2011.
- [25] M. Tepper and G. Sapiro, "Compressed nonnegative matrix factorization is fast and accurate," *IEEE Trans. Signal Process.*, vol. 64, no. 9, pp. 2269–2283, May 2016.
- [26] H. Gao, F. Nie, W. Cai, and H. Huang, "Robust capped norm nonnegative matrix factorization: Capped norm NMF," in *Proc. CIKM*, Melbourne, Australia, Oct. 2015, pp. 871–880.
- [27] J. Huang, F. Nie, H. Huang, and C. Ding, "Robust manifold nonnegative matrix factorization," *ACM Trans. Knowl. Discov. Data*, vol. 8, no. 3, pp. 1–21, Jun. 2014.
- [28] S. Yang, C. Hou, C. Zhang, and Y. Wu, "Robust non-negative matrix factorization via joint sparse and graph regularization for transfer learning," *Neural Comput. and Appl.*, vol. 23, no. 2, pp. 541–559, 2013.
- [29] L. Du, X. Li, and Y.-D. Shen, "Robust nonnegative matrix factorization via half-quadratic minimization," in *Proc. ICDM*, Brussels, Belgium, Dec. 2012, pp. 201–210.
- [30] C. Ding and D. Kong, "Nonnegative matrix factorization using a robust error function," in *Proc. ICASSP*, Kyoto, Japan, Mar. 2012, pp. 2033–2036.



Fig. 9. Results of background-foreground separation by all the online and batch algorithms on the Hall (upper) and Escalator (lower) video sequences. The labels (a) to (g) denote the original video frames and results produced by ONMF, ORPCA, OADMM, OPGD, BADMM, BPGD respectively. For each algorithm, the left column denotes the background and the right column denotes the foreground (moving objects).

- [31] D. Kong, C. Ding, and H. Huang, "Robust nonnegative matrix factorization using ℓ_{21} -norm," in *Proc. CIKM*, Glasgow, Scotland, UK, Oct. 2011, pp. 673–682.
- [32] A. Cichocki, S. Cruces, and S.-i. Amari, "Generalized alpha-beta divergences and their application to robust nonnegative matrix factorization," *Entropy*, vol. 13, no. 1, pp. 134–170, 2011.
- [33] S. P. Kasiviswanathan, H. Wang, A. Banerjee, and P. Melville, "Online ℓ_1 -dictionary learning with application to novel document detection," in *Proc. NIPS*, Lake Tahoe, USA, Dec. 2012, pp. 2258–2266.
- [34] C. Févotte and N. Dobigeon, "Nonlinear hyperspectral unmixing with robust nonnegative matrix factorization," *IEEE Trans. Image Process.*, vol. 24, no. 12, pp. 4810–4819, 2015.
- [35] B. Shen, B. Liu, Q. Wang, and R. Ji, "Robust nonnegative matrix factorization via l_1 norm regularization by multiplicative updating rules," in *Proc. ICIP*, Paris, France, Oct. 2014, pp. 5282–5286.
- [36] J. Feng, H. Xu, and S. Yan, "Online robust PCA via stochastic optimization," in *Proc. NIPS*, Lake Tahoe, USA, Dec. 2013, pp. 404–412.
- [37] J. Shen, H. Xu, and P. Li, "Online optimization for max-norm regularization," in *Proc. NIPS*, Montréal, Canada, Dec. 2014, pp. 1718–1726.
- [38] N. Wang, J. Wang, and D.-Y. Yeung, "Online robust non-negative dictionary learning for visual tracking," in *Proc. ICCV*, Sydney, Australia, Dec. 2013, pp. 657–664.
- [39] R. T. Rockafellar, *Convex analysis*. Princeton University Press, 1970.
- [40] A. W. van der Vaart, *Asymptotic Statistics*. Cambridge Press, 2000.
- [41] M. Métivier, *Semimartingales: A Course on Stochastic Processes*. de Gruyter, 1982.
- [42] P. Paatero and U. Tapper, "Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values," *Environmetrics*, vol. 5, no. 2, pp. 111–126, 1994.
- [43] P. O. Hoyer, "Non-negative matrix factorization with sparseness constraints," *J. Mach. Learn. Res.*, vol. 5, pp. 1457–1469, Dec. 2004.
- [44] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online learning for matrix factorization and sparse coding," *J. Mach. Learn. Res.*, vol. 11, pp. 19–60, Mar 2010.
- [45] D. Wang and H. Lu, "On-line learning parts-based representation via incremental orthogonal projective non-negative matrix factorization," *Signal Process.*, vol. 93, no. 6, pp. 1608 – 1623, 2013.
- [46] J. Xing, J. Gao, B. Li, W. Hu, and S. Yan, "Robust object tracking with online multi-lifespan dictionary learning," in *Proc. ICCV*, Sydney, Australia, Dec. 2013, pp. 665–672.
- [47] S. Zhang, S. Kasiviswanathan, P. C. Yuen, and M. Harandi, "Online dictionary learning on symmetric positive definite manifolds with vision applications," in *Proc. AAAI*, Austin, Texas, Jan. 2015, pp. 3165–3173.
- [48] X. Zhang, N. Guan, D. Tao, X. Qiu, and Z. Luo, "Online multi-modal robust non-negative dictionary learning for visual tracking," *PLoS ONE*, vol. 10, no. 5, pp. 1–17, 2015.
- [49] J. Zhan, B. Lois, and N. Vaswani, "Online (and offline) robust PCA: Novel algorithms and performance guarantees," arXiv:1601.07985, 2016.

- [50] B. Lois and N. Vaswani, "Online matrix completion and online robust PCA," arXiv:1503.03525, 2015.
- [51] X. Guo, "Online robust low rank matrix recovery," in *Proc. IJCAI*, Buenos Aires, Argentina, Jul. 2015, pp. 3540–3546.
- [52] P. Sprechmann, A. Bronstein, and G. Sapiro, "Real-time online singing voice separation from monaural recordings using robust low-rank modeling," in *Proc. ISMIR*, Porto, Portugal, Oct. 2012, pp. 67–72.
- [53] P. Sprechmann, A. M. Bronstein, and G. Sapiro, "Learning efficient sparse and low rank models," arXiv:1212.3631, 2012.
- [54] Y. Tsybkin, *Adaptation and Learning in automatic systems*. Academic Press, New York, 1971.
- [55] L. Bottou, "Online algorithms and stochastic approximations," in *Online Learning and Neural Networks*. Cambridge University Press, 1998.
- [56] L. Bottou and O. Bousquet, "The tradeoffs of large scale learning," in *Proc. NIPS*, Vancouver, Canada, Dec. 2008, pp. 161–168.
- [57] A. Shapiro and A. Philpott, "A tutorial on stochastic programming," <http://stoprog.org/sites/default/files/SPTutorial/TutorialSP.pdf>, 2007.
- [58] C.-J. Lin, "On the convergence of multiplicative update algorithms for nonnegative matrix factorization," *IEEE Trans. Neural Netw.*, vol. 18, no. 6, pp. 1589–1596, 2007.
- [59] A. Mensch, J. Mairal, B. Thirion, and G. Varoquaux, "Dictionary learning for massive matrix factorization," in *Proc. ICML*, New York, USA, Jul. 2016.
- [60] J. Mairal, "Stochastic majorization-minimization algorithms for large-scale optimization," in *Proc. NIPS*, Lake Tahoe, USA, Dec. 2013, pp. 2283–2291.
- [61] M. Razaviyayn, M. Sanjabi, and Z.-Q. Luo, "A stochastic successive minimization method for nonsmooth nonconvex optimization with applications to transceiver design in wireless communication networks," *Math. Program.*, vol. 157, no. 2, pp. 515–545, 2016.
- [62] D. Sun and C. Fevotte, "Alternating direction method of multipliers for non-negative matrix factorization with the beta-divergence," in *Proc. ICASSP*, Florence, Italy, May 2014, pp. 6201–6205.
- [63] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM J. Optim.*, vol. 23, no. 2, pp. 1126–1153, 2013.
- [64] M. Hong, M. Razaviyayn, Z. Q. Luo, and J. S. Pang, "A unified algorithmic framework for block-structured optimization involving big data: With applications in machine learning and signal processing," *IEEE Signal Process. Mag.*, vol. 33, no. 1, pp. 57–77, Jan. 2016.
- [65] H. Zhang and L. Cheng, "Projected shrinkage algorithm for box-constrained ℓ_1 -minimization," *Optim. Lett.*, vol. 1, no. 1, pp. 1–16, 2015.
- [66] P. Tseng, "On accelerated proximal gradient methods for convex-concave optimization," Technical report, University of Washington, Seattle, 2008.
- [67] Y. Nesterov, "Gradient methods for minimizing composite functions," *Math. Program.*, vol. 140, no. 1, pp. 125–161, 2013.
- [68] S. A. Vavasis, "On the complexity of nonnegative matrix factorization," *SIAM J. Optim.*, vol. 20, no. 3, pp. 1364–1377, 2009.
- [69] D. Hajinezhad, T. H. Chang, X. Wang, Q. Shi, and M. Hong, "Non-negative matrix factorization using admm: Algorithm and convergence analysis," in *Proc. ICASSP*, Shanghai, China, Mar. 2016, pp. 4742–4746.
- [70] D. Donoho and V. Stodden, "When does non-negative matrix factorization give correct decomposition into parts?" in *Proc. NIPS*, Vancouver, Canada, Dec. 2004, pp. 1141–1148.
- [71] J. E. Cohen and U. G. Rothblum, "Nonnegative ranks, decompositions, and factorizations of nonnegative matrices," *Linear Algebra Appl.*, vol. 190, no. 1, pp. 149–168, 1993.
- [72] S. Arora, R. Ge, R. Kannan, and A. Moitra, "Computing a nonnegative matrix factorization – provably," in *Proc. STOC*, New York, New York, USA, May 2012, pp. 145–162.
- [73] N. Gillis, "Sparse and unique nonnegative matrix factorization through data preprocessing," *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 3349–3386, Nov. 2012.
- [74] A. Vandaele, N. Gillis, F. Glineur, and D. Tuytens, "Heuristics for exact nonnegative matrix factorization," *J. Glob. Optim.*, vol. 65, no. 2, pp. 369–400, 2016.
- [75] A. Nemirovski, A. Juditsky, G. Lan, and A. Shapiro, "Robust stochastic approximation approach to stochastic programming," *SIAM J. Optimiz.*, vol. 19, no. 4, pp. 1574–1609, Jan. 2009.
- [76] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra, "Efficient projections onto the l_1 -ball for learning in high dimensions," in *Proc. ICML*, Helsinki, Finland, Jul. 2008, pp. 272–279.
- [77] W. Wang and M. A. Carreira-Perpiñán, "Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application," arXiv:1309.1541, 2013.
- [78] D. M. Witten, R. Tibshirani, and T. Hastie, "A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis," *Biostatistics*, vol. 10, no. 3, pp. 515–534, 2009.
- [79] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *J. Roy. Statist. Soc. Ser. B*, vol. 67, no. 2, pp. 301–320, 2005.
- [80] V. Sindhwani and A. Ghoting, "Large-scale distributed non-negative sparse coding and sparse dictionary learning," in *Proc. KDD*, Beijing, China, Aug. 2012, pp. 489–497.
- [81] J. Duchi, "Elastic net projections," 2009. [Online]. Available: http://stanford.edu/~jduchi/projects/proj_elastic_net.pdf
- [82] J. Kim, R. Monteiro, and H. Park, "Group sparsity in nonnegative matrix factorization," in *Proc. SDM*, Anaheim, California, USA, Apr. 2012, pp. 851–862.
- [83] X. Liu, H. Lu, and H. Gu, "Group sparse non-negative matrix factorization for multi-manifold learning," in *Proc. BMVC*, Dundee, UK, Aug. 2011, pp. 1–11.
- [84] J. Friedman, T. Hastie, and R. Tibshirani, "A note on the group lasso and a sparse group lasso," arXiv:1001.0736, 2010.
- [85] N. Simon, J. Friedman, T. Hastie, and R. Tibshirani, "A sparse-group lasso," *J. Comput. Graph. Statist.*, vol. 22, pp. 231–245, 2013.
- [86] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," *J. Roy. Statist. Soc. Ser. B*, vol. 68, no. 1, pp. 49–67, 2006.
- [87] A. S. Georgiades, P. N. Belhumeur, and D. J. Kriegman, "From few to many: illumination cone models for face recognition under variable lighting and pose," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 23, no. 6, pp. 643–660, Jun. 2001.
- [88] L. Li, W. Huang, I. Y.-H. Gu, and Q. Tian, "Statistical modeling of complex backgrounds for foreground object detection," *IEEE Trans. Image Process.*, vol. 13, no. 11, pp. 1459–1472, Nov. 2004.
- [89] V. Y. F. Tan and C. Févotte, "Automatic relevance determination in nonnegative matrix factorization with the β -divergence," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 7, pp. 1592–1605, Jul. 2013.
- [90] C. M. Bishop, "Bayesian PCA," in *Proc. NIPS*, Denver, USA, Jan. 1999, pp. 382–388.
- [91] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, Jan. 2011.
- [92] T. Veldhuizen, "Measures of image quality," http://homepages.inf.ed.ac.uk/rbf/CVonline/LOCAL_COPIES/VELDHUIZEN/node18.html, 1998.
- [93] E. Candès, X. Li, Y. Ma, and J. Wright, "Robust principal component analysis?" *J. ACM*, vol. 58, no. 3, pp. 1–37, Jun. 2011.
- [94] P. Netrapalli, N. U N, S. Sanghavi, A. Anandkumar, and P. Jain, "Non-convex robust pca," in *Proc. NIPS*, Montréal, Canada, Dec. 2014, pp. 1107–1115.
- [95] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence. With application to music analysis," *Neural Comput.*, vol. 21, no. 3, pp. 793–830, Mar. 2009.
- [96] C. Févotte and J. Idier, "Algorithms for nonnegative matrix factorization with the beta-divergence," *Neural Comput.*, vol. 23, no. 9, pp. 2421–2456, 2011.
- [97] C. Bao, H. Ji, Y. Quan, and Z. Shen, "Dictionary learning for sparse coding: Algorithms and analysis," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 7, pp. 1356–1369, July 2015.
- [98] J. F. Bonnans and A. Shapiro, "Optimization problems with perturbations: A guided tour," *SIAM Rev.*, no. 2, pp. 228–264, 1998.
- [99] K. Sydsæter, P. Hammond, A. Seierstad, and A. Strom, *Further Mathematics for Economic Analysis*. Pearson, 2005.
- [100] P. Billingsley, *Probability and Measure*, 2nd ed. John Wiley & Sons, 1986.
- [101] D. L. Fisk, "Quasi-martingales," *Trans. Amer. Math. Soc.*, vol. 120, no. 3, pp. 369–89, 1965.

Supplemental Material

In this supplemental material, the indices of all the sections, definitions, lemmas and equations are prepended with an ‘S’ to distinguish those in the main text. The organization of this article is as follows. In section S-1, we drive the ADMM algorithms presented in Section IV-B. In Section S-2, we extend our two solvers (based on PGD and ADMM) in Section IV to the *batch* NMF problems with outliers. Then, we provide detailed proofs of the theorems and lemmas in Section V in Section S-3 to S-8. The technical lemmas used in the algorithm derivation (Section IV) and the convergence analysis (Section V) are shown in Section S-9. Finally, we show additional experiment results in Section S-10 to supplement those in Section VII. The *finite* constants c , c_1 and c_2 are used repeatedly in different sections, and their meanings depend on the context.

S-1. DERIVATION OF THE ADMM ALGORITHMS IN SECTION IV-B

A. Algorithms for (9)

Minimizing \mathbf{h} and \mathbf{r} amounts to solving the following two unconstrained problems

$$\begin{aligned} & \min_{\mathbf{h}} \frac{1}{2} \|\mathbf{W}\mathbf{h} + (\mathbf{r} - \mathbf{v})\|_2^2 + \boldsymbol{\alpha}^T(\mathbf{h} - \mathbf{u}) + \frac{\rho_1}{2} \|\mathbf{h} - \mathbf{u}\|_2^2 \\ \iff & \min_{\mathbf{h}} \frac{1}{2} \mathbf{h}^T (\mathbf{W}^T \mathbf{W} + \rho_1 \mathbf{I}) \mathbf{h} + (\mathbf{W}^T (\mathbf{r} - \mathbf{v}) - \rho_1 \mathbf{u} + \boldsymbol{\alpha})^T \mathbf{h} \end{aligned} \quad (\text{S-1})$$

and

$$\begin{aligned} & \min_{\mathbf{r}} \frac{1}{2} \|(\mathbf{v} - \mathbf{W}\mathbf{h}) - \mathbf{r}\|_2^2 + \frac{\tilde{\rho}_2}{2} \|\mathbf{r} - \mathbf{q}\|_2^2 + \boldsymbol{\beta}^T(\mathbf{r} - \mathbf{q}) + \lambda \|\mathbf{r}\|_1 \\ \iff & \min_{\mathbf{r}} \frac{1 + \tilde{\rho}_2}{2} \mathbf{r}^T \mathbf{r} + (\boldsymbol{\beta} - \mathbf{v} + \mathbf{W}\mathbf{h} - \tilde{\rho}_2 \mathbf{q})^T \mathbf{r} + \lambda \|\mathbf{r}\|_1 \\ \iff & \min_{\mathbf{r}} \frac{1 + \tilde{\rho}_2}{2} \left\| \mathbf{r} - \frac{\tilde{\rho}_2 \mathbf{q} + \mathbf{v} - \boldsymbol{\beta} - \mathbf{W}\mathbf{h}}{1 + \tilde{\rho}_2} \right\|_2^2 + \lambda \|\mathbf{r}\|_1. \end{aligned} \quad (\text{S-2})$$

We notice that (S-1) is a standard strongly convex quadratic minimization problem and (S-2) is a standard proximal minimization problem with ℓ_1 norm, thus the closed-form optimal solutions for (S-1) and (S-2) are

$$\begin{aligned} \mathbf{h}^* &= (\mathbf{W}^T \mathbf{W} + \rho_1 \mathbf{I})^{-1} (\mathbf{W}^T (\mathbf{v} - \mathbf{r}) + \rho_1 \mathbf{u} - \boldsymbol{\alpha}) \\ \mathbf{r}^* &= \mathcal{S}_{\lambda/(1+\tilde{\rho}_2)} \left(\frac{\tilde{\rho}_2 \mathbf{q} + \mathbf{v} - \boldsymbol{\beta} - \mathbf{W}\mathbf{h}}{1 + \tilde{\rho}_2} \right) = \frac{\mathcal{S}_{\lambda}(\tilde{\rho}_2 \mathbf{q} + \mathbf{v} - \boldsymbol{\beta} - \mathbf{W}\mathbf{h})}{1 + \tilde{\rho}_2}. \end{aligned}$$

Minimizing \mathbf{u} and \mathbf{q} amounts to solving the following two constrained problems

$$\min_{\mathbf{u} \geq 0} \frac{\rho_1}{2} \|\mathbf{h} - \mathbf{u}\|_2^2 - \boldsymbol{\alpha}^T \mathbf{u} \quad (\text{S-3})$$

$$\min_{\|\mathbf{q}\|_{\infty} \leq M} \frac{\tilde{\rho}_2}{2} \|\mathbf{r} - \mathbf{q}\|_2^2 - \boldsymbol{\beta}^T \mathbf{q}. \quad (\text{S-4})$$

Since both constraints $\mathbf{u} \geq 0$ and $\|\mathbf{q}\|_{\infty} \leq M$ are separable across coordinates, we can simply solve the unconstrained quadratic minimization problems and then project the optimal solutions to the feasible sets.

B. Algorithms for (12)

Minimizing \mathbf{W} amounts to solving the unconstrained quadratic minimization problem

$$\begin{aligned} & \min_{\mathbf{W}} \frac{1}{2} \text{tr}(\mathbf{W}^T \mathbf{W} \mathbf{A}_t) - \text{tr}(\mathbf{W}^T \mathbf{B}_t) + \langle \mathbf{D}, \mathbf{W} - \mathbf{Q} \rangle + \frac{\tilde{\rho}_3}{2} \|\mathbf{W} - \mathbf{Q}\|_F^2 \\ \iff & \min_{\mathbf{W}} \frac{1}{2} \text{tr}(\mathbf{W} (\mathbf{A}_t + \tilde{\rho}_3 \mathbf{I}) \mathbf{W}^T) - \text{tr}((\mathbf{B}_t - \mathbf{D} + \tilde{\rho}_3 \mathbf{Q})^T \mathbf{W}) \end{aligned}$$

so

$$\mathbf{W}^* = (\mathbf{B}_t - \mathbf{D} + \tilde{\rho}_3 \mathbf{Q}) (\mathbf{A}_t + \tilde{\rho}_3 \mathbf{I})^{-1}. \quad (\text{S-5})$$

Minimizing \mathbf{Q} amounts to solving the constrained quadratic minimization problem

$$\begin{aligned} & \min_{\mathbf{Q} \in \mathcal{C}} \frac{\tilde{\rho}_3}{2} \|\mathbf{W} - \mathbf{Q}\|_F^2 - \langle \mathbf{D}, \mathbf{Q} \rangle \\ \iff & \min_{\mathbf{Q} \in \mathcal{C}} \frac{\tilde{\rho}_3}{2} \|\mathbf{Q} - (\mathbf{W} + \mathbf{D}/\tilde{\rho}_3)\|_F^2. \end{aligned}$$

Then we have

$$\mathbf{Q}^* = \mathcal{P}_{\mathcal{C}}(\mathbf{W} + \mathbf{D}/\tilde{\rho}_3). \quad (\text{S-6})$$

S-2. EXTENSION TO THE BATCH NMF PROBLEM WITH OUTLIERS

A. Problem Formulation

As usual, we denote the data matrix (with outliers) as \mathbf{V} , basis matrix as \mathbf{W} , coefficient matrix as \mathbf{H} and outlier matrix as \mathbf{R} . The number of total data samples is denoted as N . Then, the batch counterpart for the online NMF problem with outliers can be formulated as follows

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{V} - \mathbf{W}\mathbf{H} - \mathbf{R}\|_F^2 + \lambda \|\mathbf{R}\|_{1,1} \\ \text{s. t.} \quad & \mathbf{H} \in \mathbb{R}_+^{K \times N}, \mathbf{R} \in \tilde{\mathcal{R}}, \mathbf{W} \in \mathcal{C}, \end{aligned} \quad (\text{S-7})$$

where $\tilde{\mathcal{R}} = \{\mathbf{R} \in \mathbb{R}^{F \times N} \mid |r_{i,j}| \leq M, \forall (i,j) \in [F] \times [N]\}$.

B. Notations

In the sequel, we overload soft-thresholding operators $\tilde{\mathcal{S}}_{\lambda,M}$ and \mathcal{S}_λ . When these two operators are applied to matrices, each operator denotes entrywise soft-thresholding. The updated variables are denoted with superscripts ‘+’.

C. Batch algorithm based on PGD (BPGD)

Based on the principle of block coordinate descent, we update \mathbf{H} , \mathbf{R} and \mathbf{W} sequentially as follows.

$$\mathbf{H}^+ := \mathcal{P}_+(\mathbf{H} - \eta_1(\mathbf{W})\mathbf{W}^T(\mathbf{W}\mathbf{H} + \mathbf{R} - \mathbf{V})) \quad (\text{S-8})$$

$$\mathbf{R}^+ := \tilde{\mathcal{S}}_{\lambda,M}(\mathbf{V} - \mathbf{W}\mathbf{H}^+) \quad (\text{S-9})$$

$$\mathbf{W}_{:j}^+ := \frac{\mathcal{P}_+(\mathbf{W} - \eta_2(\mathbf{H}^+)\mathbf{G})_{:j}}{\max\{1, \|\mathcal{P}_+(\mathbf{W} - \eta_2(\mathbf{H}^+)\mathbf{G})_{:j}\|_2\}}, \quad \forall j \in [K], \quad (\text{S-10})$$

where $\mathbf{G} = (\mathbf{W}\mathbf{H}^+ + \mathbf{R}^+ - \mathbf{V})\mathbf{H}^{+T}$, $\eta_1(\mathbf{W}) \in (0, \|\mathbf{W}\|_2^{-2}]$ and $\eta_2(\mathbf{H}^+) \in (0, \|\mathbf{H}^+\mathbf{H}^{+T}\|_F^{-1}]$.

D. Batch algorithm based on ADMM (BADMM)

(S-7) can be reformulated as

$$\begin{aligned} \min \quad & \frac{1}{2} \|\mathbf{V} - \mathbf{W}\mathbf{H} - \mathbf{R}\|_F^2 + \lambda \|\mathbf{R}\|_{1,1} \\ \text{s. t.} \quad & \mathbf{H} = \mathbf{U}, \mathbf{R} = \mathbf{Q}, \mathbf{W} = \mathbf{\Psi}, \mathbf{U} \in \mathbb{R}_+^{K \times N}, \mathbf{Q} \in \tilde{\mathcal{R}}, \mathbf{\Psi} \in \mathcal{C}. \end{aligned}$$

Thus the augmented Lagrangian is

$$\begin{aligned} \tilde{\mathcal{L}}(\mathbf{H}, \mathbf{R}, \mathbf{W}, \mathbf{U}, \mathbf{Q}, \mathbf{\Psi}, \mathbf{A}, \mathbf{B}, \mathbf{D}) = & \frac{1}{2} \|\mathbf{V} - \mathbf{W}\mathbf{H} - \mathbf{R}\|_F^2 + \lambda \|\mathbf{R}\|_{1,1} + \langle \mathbf{A}, \mathbf{H} - \mathbf{U} \rangle + \langle \mathbf{B}, \mathbf{R} - \mathbf{Q} \rangle + \langle \mathbf{D}, \mathbf{W} - \mathbf{\Psi} \rangle \\ & + \frac{\tilde{\rho}_1}{2} \|\mathbf{H} - \mathbf{U}\|_F^2 + \frac{\tilde{\rho}_2}{2} \|\mathbf{R} - \mathbf{Q}\|_F^2 + \frac{\tilde{\rho}_3}{2} \|\mathbf{W} - \mathbf{\Psi}\|_F^2, \end{aligned} \quad (\text{S-11})$$

where \mathbf{A} , \mathbf{B} and \mathbf{D} are dual variables and $\tilde{\rho}_1$, $\tilde{\rho}_2$ and $\tilde{\rho}_3$ are positive penalty parameters. Therefore we can derive the following update rules

$$\mathbf{H}^+ := (\mathbf{W}^T\mathbf{W} + \tilde{\rho}_1\mathbf{I})^{-1} (\mathbf{W}^T(\mathbf{V} - \mathbf{R}) + \tilde{\rho}_1\mathbf{U} - \mathbf{A}) \quad (\text{S-12})$$

$$\mathbf{R}^+ := \mathcal{S}_\lambda(\tilde{\rho}_2\mathbf{Q} + \mathbf{V} - \mathbf{B} - \mathbf{W}\mathbf{H}^+)/ (1 + \tilde{\rho}_2) \quad (\text{S-13})$$

$$\mathbf{W}^+ := \left((\mathbf{V} - \mathbf{R}^+)\mathbf{H}^{+T} - \mathbf{D} + \tilde{\rho}_3\mathbf{\Psi} \right) \left(\mathbf{H}^+\mathbf{H}^{+T} + \tilde{\rho}_3\mathbf{I} \right)^{-1} \quad (\text{S-14})$$

$$\mathbf{U}^+ := \mathcal{P}_+(\mathbf{H}^+ + \mathbf{A}/\tilde{\rho}_1) \quad (\text{S-15})$$

$$\mathbf{Q}^+ := \mathcal{P}_{\tilde{\mathcal{R}}}(\mathbf{R}^+ + \mathbf{B}/\tilde{\rho}_2) \quad (\text{S-16})$$

$$\mathbf{\Psi}^+ := \mathcal{P}_{\mathcal{C}}(\mathbf{W}^+ + \mathbf{D}/\tilde{\rho}_3) \quad (\text{S-17})$$

$$\mathbf{A}^+ := \mathbf{A} + \tilde{\rho}_1(\mathbf{H}^+ - \mathbf{U}^+) \quad (\text{S-18})$$

$$\mathbf{B}^+ := \mathbf{B} + \tilde{\rho}_2(\mathbf{R}^+ - \mathbf{Q}^+) \quad (\text{S-19})$$

$$\mathbf{D}^+ := \mathbf{D} + \tilde{\rho}_3(\mathbf{W}^+ - \mathbf{\Psi}^+), \quad (\text{S-20})$$

S-3. PROOF OF LEMMA 1

Before proving Lemma 1, we present two lemmas which will be used in the proof. Both lemmas can be proved using straightforward calculations. See Section S-7 and S-8 for detailed proofs.

Lemma S-1. *If for each $\mathbf{v} \in \mathcal{V}$, both $\mathbf{h}^*(\mathbf{v}, \mathbf{W})$ and $\mathbf{r}^*(\mathbf{v}, \mathbf{W})$ in (S-21) are Lipschitz on \mathcal{C} , with Lipschitz constants c_1 and c_2 (independent of \mathbf{v}) respectively, then $\mathbf{W} \mapsto \nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$ is Lipschitz on \mathcal{C} with Lipschitz constant c_3 (independent of \mathbf{v}). Consequently, $\nabla f(\mathbf{W})$ in (S-22) is Lipschitz on \mathcal{C} with Lipschitz constant c_3 .*

Lemma S-2. *Let $\mathbf{z}, \mathbf{z}' \in \mathcal{Z} \subseteq \mathbb{R}^m$ and $\mathbf{A}, \mathbf{A}' \in \mathcal{A} \subseteq \mathbb{R}^{m \times n}$, where both \mathcal{Z} and \mathcal{A} are compact sets. Let \mathcal{B} be a compact set in \mathbb{R}^n , and define $g : \mathcal{B} \rightarrow \mathbb{R}$ as $g(\mathbf{b}) = 1/2 \|\mathbf{z} - \mathbf{A}\mathbf{b}\|_2^2 - 1/2 \|\mathbf{z}' - \mathbf{A}'\mathbf{b}\|_2^2$. Then g is Lipschitz on \mathcal{B} with Lipschitz constant $c_1 \|\mathbf{z} - \mathbf{z}'\|_2 + c_2 \|\mathbf{A} - \mathbf{A}'\|_2$, where c_1 and c_2 are two positive constants. In particular, when both \mathbf{z}' and \mathbf{A}' are zero, we have that $\tilde{g}(\mathbf{b}) = 1/2 \|\mathbf{z} - \mathbf{A}\mathbf{b}\|_2^2$ is Lipschitz on \mathcal{B} with Lipschitz constant c independent of \mathbf{z} and \mathbf{A} .*

It is easy to verify that the following conditions hold

- 1) $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$ is differentiable on $\mathcal{V} \times \mathcal{C}$, for each $(\mathbf{h}, \mathbf{r}) \in \mathcal{H} \times \mathcal{R}$,
- 2) $\tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$ and $\nabla_{(\mathbf{v}, \mathbf{W})} \tilde{\ell}(\mathbf{v}, \mathbf{W}, \mathbf{h}, \mathbf{r})$ are continuous on $\mathcal{V} \times \mathcal{C} \times \mathcal{H} \times \mathcal{R}$,
- 3) (38) has unique minimizer $(\mathbf{h}^*(\mathbf{v}, \mathbf{W}), \mathbf{r}^*(\mathbf{v}, \mathbf{W}))$ for each $(\mathbf{v}, \mathbf{W}) \in \mathcal{V} \times \mathcal{C}$, due to Assumption 2.

Thus, we can invoke Danskin's theorem (see Lemma S-5) to conclude that $\ell(\mathbf{v}, \mathbf{W})$ is differentiable on $\mathcal{V} \times \mathcal{C}$ and

$$\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}) = (\mathbf{W}\mathbf{h}^*(\mathbf{v}, \mathbf{W}) + \mathbf{r}^*(\mathbf{v}, \mathbf{W}) - \mathbf{v}) \mathbf{h}^*(\mathbf{v}, \mathbf{W})^T. \quad (\text{S-21})$$

Furthermore, we can show $(\mathbf{h}^*(\mathbf{v}, \mathbf{W}), \mathbf{r}^*(\mathbf{v}, \mathbf{W}))$ is continuous on $\mathcal{V} \times \mathcal{C}$ by the maximum theorem (see Lemma S-6), since the conditions in this theorem are trivially satisfied in our case. Thus, $\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$ is continuous on $\mathcal{V} \times \mathcal{C}$.

Leveraging the regularity of $\ell(\mathbf{v}, \mathbf{W})$, we proceed to show the regularity of $f(\mathbf{W})$. Since for all $\mathbf{v} \in \mathcal{V}$, both $\ell(\mathbf{v}, \mathbf{W})$ and $\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$ are continuous on \mathcal{C} , by Leibniz integral rule (see Lemma S-7), we conclude that $f(\mathbf{W})$ is differentiable on \mathcal{C} and

$$\nabla f(\mathbf{W}) = \mathbb{E}_{\mathbf{v}}[\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})]. \quad (\text{S-22})$$

By Lemma S-1, to show both $\mathbf{W} \mapsto \nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W})$ and $\nabla f(\mathbf{W})$ are Lipschitz on \mathcal{C} , it suffices to show both $\mathbf{h}^*(\mathbf{v}, \mathbf{W})$ and $\mathbf{r}^*(\mathbf{v}, \mathbf{W})$ are Lipschitz on \mathcal{C} , for all $\mathbf{v} \in \mathcal{V}$. Fix arbitrary $\mathbf{v}_1, \mathbf{v}_2 \in \mathcal{V}$ and $\mathbf{W}_1, \mathbf{W}_2 \in \mathcal{C}$. Define

$$\begin{aligned} d(\mathbf{h}, \mathbf{r}) &\triangleq \tilde{\ell}(\mathbf{v}_1, \mathbf{W}_1, \mathbf{h}, \mathbf{r}) - \tilde{\ell}(\mathbf{v}_2, \mathbf{W}_2, \mathbf{h}, \mathbf{r}) \\ &= \frac{1}{2} \|\mathbf{v}_1 - \mathbf{Y}_1 \mathbf{b}(\mathbf{h}, \mathbf{r})\|_2^2 - \frac{1}{2} \|\mathbf{v}_2 - \mathbf{Y}_2 \mathbf{b}(\mathbf{h}, \mathbf{r})\|_2^2 \end{aligned}$$

where $\mathbf{Y}_i = [\mathbf{W}_i \ \mathbf{I}]$, $i = 1, 2$ and $\mathbf{b}(\mathbf{h}, \mathbf{r}) = [\mathbf{h}^T \ \mathbf{r}^T]^T$. By Lemma S-2, we have for all $(\mathbf{h}_1, \mathbf{r}_1), (\mathbf{h}_2, \mathbf{r}_2) \in \mathcal{H} \times \mathcal{R}$,

$$|d(\mathbf{h}_1, \mathbf{r}_1) - d(\mathbf{h}_2, \mathbf{r}_2)| \leq (c_1 \|\mathbf{v}_1 - \mathbf{v}_2\|_2 + c_2 \|\mathbf{Y}_1 - \mathbf{Y}_2\|_2) \|\mathbf{b}(\mathbf{h}_1, \mathbf{r}_1) - \mathbf{b}(\mathbf{h}_2, \mathbf{r}_2)\|_2 \quad (\text{S-23})$$

where c_1 and c_2 are positive constants. In particular, we have

$$|d(\mathbf{h}_1^*, \mathbf{r}_1^*) - d(\mathbf{h}_2^*, \mathbf{r}_2^*)| \leq (c_1 \|\mathbf{v}_1 - \mathbf{v}_2\|_2 + c_2 \|\mathbf{Y}_1 - \mathbf{Y}_2\|_2) \|\mathbf{b}(\mathbf{h}_1^*, \mathbf{r}_1^*) - \mathbf{b}(\mathbf{h}_2^*, \mathbf{r}_2^*)\|_2 \quad (\text{S-24})$$

where $\mathbf{h}_i^* = \mathbf{h}^*(\mathbf{v}_i, \mathbf{W}_i)$ and $\mathbf{r}_i^* = \mathbf{r}^*(\mathbf{v}_i, \mathbf{W}_i)$, $i = 1, 2$. On the other hand, by Assumption 2,

$$\begin{aligned} |d(\mathbf{h}_2^*, \mathbf{r}_2^*) - d(\mathbf{h}_1^*, \mathbf{r}_1^*)| &= \left| \tilde{\ell}(\mathbf{v}_1, \mathbf{W}_1, \mathbf{h}_2^*, \mathbf{r}_2^*) - \tilde{\ell}(\mathbf{v}_2, \mathbf{W}_2, \mathbf{h}_2^*, \mathbf{r}_2^*) - \tilde{\ell}(\mathbf{v}_1, \mathbf{W}_1, \mathbf{h}_1^*, \mathbf{r}_1^*) + \tilde{\ell}(\mathbf{v}_2, \mathbf{W}_2, \mathbf{h}_1^*, \mathbf{r}_1^*) \right| \\ &= (\tilde{\ell}(\mathbf{v}_1, \mathbf{W}_1, \mathbf{h}_2^*, \mathbf{r}_2^*) - \tilde{\ell}(\mathbf{v}_1, \mathbf{W}_1, \mathbf{h}_1^*, \mathbf{r}_1^*)) + (\tilde{\ell}(\mathbf{v}_2, \mathbf{W}_2, \mathbf{h}_1^*, \mathbf{r}_1^*) - \tilde{\ell}(\mathbf{v}_2, \mathbf{W}_2, \mathbf{h}_2^*, \mathbf{r}_2^*)) \\ &\geq m_1 \|\mathbf{b}(\mathbf{h}_1^*, \mathbf{r}_1^*) - \mathbf{b}(\mathbf{h}_2^*, \mathbf{r}_2^*)\|_2^2. \end{aligned} \quad (\text{S-25})$$

Combining (S-24) and (S-25), we have

$$\max\{\|\mathbf{h}_1^* - \mathbf{h}_2^*\|_2, \|\mathbf{r}_1^* - \mathbf{r}_2^*\|_2\} \leq \|\mathbf{b}(\mathbf{h}_1^*, \mathbf{r}_1^*) - \mathbf{b}(\mathbf{h}_2^*, \mathbf{r}_2^*)\|_2 \leq c'_1 \|\mathbf{v}_1 - \mathbf{v}_2\|_2 + c'_2 \|\mathbf{W}_1 - \mathbf{W}_2\|_2 \quad (\text{S-26})$$

where $c'_i = c_i/m_1$, $i = 1, 2$. This indeed shows both $\mathbf{h}^*(\mathbf{v}, \mathbf{W})$ and $\mathbf{r}^*(\mathbf{v}, \mathbf{W})$ are Lipschitz on $\mathcal{V} \times \mathcal{C}$ since

$$c'_1 \|\mathbf{v}_1 - \mathbf{v}_2\|_2 + c'_2 \|\mathbf{W}_1 - \mathbf{W}_2\|_2 \leq 2 \max(c'_1, c'_2) \|[\mathbf{v}_1 \ \mathbf{W}_1] - [\mathbf{v}_2 \ \mathbf{W}_2]\|_F. \quad (\text{S-27})$$

Hence we complete the proof.

S-4. PROOF OF LEMMA 2

By Assumption 3, for all $t \geq 1$, we have

$$\tilde{f}_t(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t) \geq \frac{m_2}{2} \|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F^2, \quad (\text{S-28})$$

since $\mathbf{W}_t = \arg \min_{\mathbf{W}} \tilde{f}_t(\mathbf{W})$. On the other hand,

$$\begin{aligned} \tilde{f}_t(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t) &= \tilde{f}_t(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_{t+1}) + \tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t) + \tilde{f}_{t+1}(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \\ &\leq \tilde{d}_t(\mathbf{W}_{t+1}) - \tilde{d}_t(\mathbf{W}_t) \end{aligned}$$

where $\tilde{d}_t(\mathbf{W}) \triangleq \tilde{f}_t(\mathbf{W}) - \tilde{f}_{t+1}(\mathbf{W})$, for all $\mathbf{W} \in \mathcal{C}$. We aim to show \tilde{d}_t is Lipschitz on \mathcal{C} with Lipschitz constant only dependent on t . For all $\mathbf{W} \in \mathcal{C}$, we have

$$\begin{aligned} \tilde{d}_t(\mathbf{W}) &= \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 + \lambda \|\mathbf{r}_i\|_1 - \frac{1}{t+1} \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 + \lambda \|\mathbf{r}_i\|_1 \\ &= \frac{1}{t(t+1)} \left((t+1) \sum_{i=1}^t \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 + \lambda \|\mathbf{r}_i\|_1 - t \sum_{i=1}^{t+1} \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 + \lambda \|\mathbf{r}_i\|_1 \right) \\ &\stackrel{\text{c}}{=} \frac{1}{t(t+1)} \left(\sum_{i=1}^t \frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 - \frac{t}{2} \|\mathbf{v}_{t+1} - \mathbf{W}\mathbf{h}_{t+1} - \mathbf{r}_{t+1}\|_2^2 \right) \\ &= \frac{1}{t(t+1)} \sum_{i=1}^t \left(\frac{1}{2} \|\mathbf{v}_i - \mathbf{W}\mathbf{h}_i - \mathbf{r}_i\|_2^2 - \frac{1}{2} \|\mathbf{v}_{t+1} - \mathbf{W}\mathbf{h}_{t+1} - \mathbf{r}_{t+1}\|_2^2 \right), \end{aligned}$$

where $\stackrel{\text{c}}{=}$ denotes equality up to an additive constant (independent of \mathbf{W}). By a reasoning similar to the one for Lemma S-2, we can show there exist some positive constants c_1 and c_2 such that

$$\begin{aligned} \left| \tilde{d}_t(\mathbf{W}_{t+1}) - \tilde{d}_t(\mathbf{W}_t) \right| &\leq \frac{1}{t(t+1)} \sum_{i=1}^t (c_1 \|\mathbf{h}_i - \mathbf{h}_{t+1}\|_2 + c_2 \|\mathbf{v}_i - \mathbf{v}_{t+1}\|_2 + c_2 \|\mathbf{r}_i - \mathbf{r}_{t+1}\|_2) \|\mathbf{W}_{t+1} - \mathbf{W}_t\|_2 \\ &\leq \frac{c_3}{t+1} \|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F, \end{aligned} \tag{S-29}$$

where $c_3 > 0$ is a constant since for all $i \geq 1$, \mathbf{v}_i , \mathbf{h}_i and \mathbf{r}_i are bounded a.s.. Combining (S-28) and (S-29), we have with probability one,

$$\|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F \leq \frac{c'_3}{t+1}, \tag{S-30}$$

where $c'_3 = c_3/m_2$. Hence we complete the proof.

S-5. PROOF OF THEOREM 1

We prove that $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$ converges a.s. by the quasi-martingale convergence theorem (see Lemma S-8). Let us first define a filtration $\{\mathcal{F}_t\}_{t \geq 1}$ where $\mathcal{F}_t \triangleq \sigma\{\mathbf{v}_i, \mathbf{W}_i, \mathbf{h}_i, \mathbf{r}_i\}_{i \in [t]}$ is the minimal σ -algebra such that $\{\mathbf{v}_i, \mathbf{W}_i, \mathbf{h}_i, \mathbf{r}_i\}_{i \in [t]}$ are measurable. Also define

$$u_t \triangleq \tilde{f}_t(\mathbf{W}_t) \quad \text{and} \quad \delta_t \triangleq \begin{cases} 1, & \text{if } \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] > 0 \\ 0, & \text{otherwise} \end{cases}, \tag{S-31}$$

then it is easy to see $\{u_t\}_{t \geq 1}$ is adapted to $\{\mathcal{F}_t\}_{t \geq 1}$. According to the quasi-martingale convergence theorem, it suffices to show $\sum_{t=1}^{\infty} \mathbb{E}[\delta_t(u_{t+1} - u_t)] < \infty$. To bound $\mathbb{E}[\delta_t(u_{t+1} - u_t)]$, we decompose $u_{t+1} - u_t$ as

$$\begin{aligned} u_{t+1} - u_t &= \tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t) + \tilde{f}_{t+1}(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \\ &= \tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t) + \frac{1}{t+1} \tilde{\ell}(\mathbf{v}_{t+1}, \mathbf{h}_{t+1}, \mathbf{r}_{t+1}, \mathbf{W}_t) + \frac{t}{t+1} \tilde{f}_t(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \\ &= \tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t) + \frac{\ell(\mathbf{v}_{t+1}, \mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t)}{t+1} \end{aligned} \tag{S-32}$$

$$= \underbrace{\tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t)}_{\langle 1 \rangle} + \frac{\ell(\mathbf{v}_{t+1}, \mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1} + \underbrace{\frac{f_t(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t)}{t+1}}_{\langle 2 \rangle}. \tag{S-33}$$

By definition, it is easy to see both $\langle 1 \rangle, \langle 2 \rangle \leq 0$. In (S-33), we insert the term $f_t(\mathbf{W}_t)$ in order to invoke Donsker's theorem (see Lemma S-9). Thus,

$$\mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] \leq \frac{\mathbb{E}[\ell(\mathbf{v}_{t+1}, \mathbf{W}_t) - f_t(\mathbf{W}_t) | \mathcal{F}_t]}{t+1} = \frac{f(\mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1}$$

and using the definition of δ_t ,

$$\mathbb{E}[\delta_t \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]] \leq \frac{\mathbb{E}[\|f(\mathbf{W}_t) - f_t(\mathbf{W}_t)\|]}{t+1} \leq \frac{\mathbb{E}[\|f - f_t\|_{\mathcal{C}}]}{t+1} = \frac{\mathbb{E}[\|\sqrt{t}(f - f_t)\|_{\mathcal{C}}]}{\sqrt{t}(t+1)}, \quad (\text{S-34})$$

where $\|\cdot\|_{\mathcal{C}}$ denotes the uniform norm on \mathcal{C} . By Lemma 1, we know for all $\mathbf{v} \in \mathcal{V}$, $\ell(\mathbf{v}, \cdot)$ is Lipschitz on \mathcal{C} with Lipschitz constant independent of \mathbf{v} . Thus, by Lemma S-11, the measurable function class $\{\ell(\cdot, \mathbf{W}) : \mathbf{W} \in \mathcal{C}\}$ is \mathbb{P} -Donsker. Consequently $\mathbb{E}[\|\sqrt{t}(f - f_t)\|_{\mathcal{C}}]$ is bounded by a constant $c > 0$. Thus, $\mathbb{E}[\delta_t \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]] \leq c/t^{3/2}$. Since $\mathbb{E}[\delta_t(u_{t+1} - u_t)] = \mathbb{E}[\mathbb{E}[\delta_t(u_{t+1} - u_t) | \mathcal{F}_t]] = \mathbb{E}[\delta_t \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]]$, we have $\sum_{t=1}^{\infty} \mathbb{E}[\delta_t(u_{t+1} - u_t)] < \infty$. Thus $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$ converges a.s.. Moreover, by Lemma S-8, we also have $\sum_{t=1}^{\infty} \mathbb{E}[\|\mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]\|] < \infty$.

Leveraging this result, we proceed to show the almost sure convergence of $\{f(\mathbf{W}_t)\}_{t \geq 1}$. By Lemma S-11, the measurable function class $\{\ell(\cdot, \mathbf{W}) : \mathbf{W} \in \mathcal{C}\}$ is also \mathbb{P} -Glivenko-Cantelli. Thus by Glivenko-Cantelli theorem (see Lemma S-10), we have $\|f_t - f\|_{\mathcal{C}} \xrightarrow{\text{a.s.}} 0$. Hence it suffices to show $\{f_t(\mathbf{W}_t)\}_{t \geq 1}$ converges a.s.. We show this by proving $f_t(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \xrightarrow{\text{a.s.}} 0$. First, from (S-33), we have

$$\begin{aligned} \frac{\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1} &= \mathbb{E}\left[\frac{\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1} \middle| \mathcal{F}_t\right] \\ &= \mathbb{E}[\tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t) | \mathcal{F}_t] + \mathbb{E}\left[\frac{\ell(\mathbf{v}_{t+1}, \mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1} \middle| \mathcal{F}_t\right] - \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] \\ &\leq \frac{\mathbb{E}[\ell(\mathbf{v}_{t+1}, \mathbf{W}_t) | \mathcal{F}_t] - f_t(\mathbf{W}_t)}{t+1} - \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] \\ &\leq \frac{f(\mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1} - \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] \\ &\leq \frac{\|f - f_t\|_{\mathcal{C}}}{t+1} - \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]. \end{aligned}$$

Since both $\sum_{t=1}^{\infty} \frac{\|f - f_t\|_{\mathcal{C}}}{t+1}$ and $\sum_{t=1}^{\infty} |\mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]|$ converge a.s., we conclude $\sum_{t=1}^{\infty} \frac{\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t)}{t+1}$ converges a.s.. Define $b_t \triangleq \tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t)$, we show $|b_{t+1} - b_t| = O(1/t)$ a.s. by proving both $|\tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t)| = O(1/t)$ and $|f_{t+1}(\mathbf{W}_{t+1}) - f_t(\mathbf{W}_t)| = O(1/t)$ a.s.. First, from (S-32) and the Lipschitz continuity of \tilde{f}_t on \mathcal{C} ,

$$\begin{aligned} |\tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t)| &\leq |\tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_{t+1}(\mathbf{W}_t)| + \left| \frac{\ell(\mathbf{v}_{t+1}, \mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t)}{t+1} \right| \\ &\leq c \|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F + \frac{|\ell(\mathbf{v}_{t+1}, \mathbf{W}_t)| + |\tilde{f}_t(\mathbf{W}_t)|}{t+1}, \end{aligned}$$

where $c > 0$ is a constant independent of t . Since both $|\ell(\mathbf{v}_{t+1}, \mathbf{W}_t)|$ and $|\tilde{f}_t(\mathbf{W}_t)|$ are bounded on \mathcal{C} a.s. and $\|\mathbf{W}_{t+1} - \mathbf{W}_t\|_F = O(1/t)$ a.s. (by Lemma 2), we have $|\tilde{f}_{t+1}(\mathbf{W}_{t+1}) - \tilde{f}_t(\mathbf{W}_t)| = O(1/t)$ a.s.. Similarly, by the Lipschitz continuity of f_t on \mathcal{C} , we also have $|f_{t+1}(\mathbf{W}_{t+1}) - f_t(\mathbf{W}_t)| = O(1/t)$ a.s.. Now we invoke Lemma S-12 to conclude

$$f_t(\mathbf{W}_t) - \tilde{f}_t(\mathbf{W}_t) \xrightarrow{\text{a.s.}} 0. \quad (\text{S-35})$$

Since $f_t(\mathbf{W}_t) - f(\mathbf{W}_t) \xrightarrow{\text{a.s.}} 0$, both $\{f(\mathbf{W}_t)\}_{t \geq 1}$ and $\{\tilde{f}_t(\mathbf{W}_t)\}_{t \geq 1}$ converge to the same almost sure limit.

S-6. PROOF OF THEOREM 2

By (S-35), it suffices to show that for every realization of $\{\mathbf{v}_t\}_{t \geq 1}$ such that $\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t) \rightarrow 0$, each subsequential limit of $\{\mathbf{W}_t\}_{t \geq 1}$ is a stationary point of f . We focus on such a realization, then all the variables in the sequel become deterministic. (With a slight abuse of notations we use the same notations to denote the deterministic variables.) By the compactness of \mathcal{V} , \mathcal{H} and \mathcal{R} , both sequences $\{\mathbf{A}_t\}_{t \geq 1}$ and $\{\mathbf{B}_t\}_{t \geq 1}$ are bounded. Thus there exist compact sets \mathcal{A} and \mathcal{B} such that $\{\mathbf{A}_t\}_{t \geq 1} \subseteq \mathcal{A}$ and $\{\mathbf{B}_t\}_{t \geq 1} \subseteq \mathcal{B}$. Similar reasoning shows that the sequence $\{\tilde{f}_t(\mathbf{0})\}_{t \geq 1}$ resides in a compact set $\tilde{\mathcal{F}}$. By the compactness of \mathcal{C} , there exists a convergent subsequence $\{\mathbf{W}_{t_m}\}_{m \geq 1}$ in $\{\mathbf{W}_t\}_{t \geq 1}$. Also, by the compactness of \mathcal{A} , \mathcal{B} and $\tilde{\mathcal{F}}$, it is possible to find convergent subsequences $\{\mathbf{A}_{t_k}\}_{k \geq 1}$, $\{\mathbf{B}_{t_k}\}_{k \geq 1}$ and $\{\tilde{f}_{t_k}(\mathbf{0})\}_{k \geq 1}$ such that $\{t_k\}_{k \geq 1} \subseteq \{t_m\}_{m \geq 1}$. Thus, we focus on the convergent sequences $\{\mathbf{W}_{t_k}\}_{k \geq 1}$, $\{\mathbf{A}_{t_k}\}_{k \geq 1}$, $\{\mathbf{B}_{t_k}\}_{k \geq 1}$ and $\{\tilde{f}_{t_k}(\mathbf{0})\}_{k \geq 1}$ and drop the subscript k to make notations uncluttered. We denote the limits of the sequences $\{\mathbf{W}_t\}_{t \geq 1}$, $\{\mathbf{A}_t\}_{t \geq 1}$ and $\{\mathbf{B}_t\}_{t \geq 1}$ as $\bar{\mathbf{W}}$, $\bar{\mathbf{A}}$ and $\bar{\mathbf{B}}$ respectively.

First, we show the sequence of differentiable functions $\{\tilde{f}_t\}_{t \geq 1}$ converges uniformly to a differentiable function \bar{f} . Since $\{\tilde{f}_t(\mathbf{0})\}_{t \geq 1}$ converges, it suffices to show the sequence $\{\nabla \tilde{f}_t\}_{t \geq 1}$ converges uniformly to a function \bar{h} . Since $\nabla \tilde{f}_t(\mathbf{W}) =$

$\mathbf{W}\mathbf{A}_t - \mathbf{B}_t$, for any $t, t' \geq 1$ and any $\mathbf{W} \in \mathcal{C}$, we have

$$\begin{aligned} \left\| \nabla \tilde{f}_t(\mathbf{W}) - \nabla \tilde{f}_{t'}(\mathbf{W}) \right\|_F &= \left\| \mathbf{W}(\mathbf{A}_t - \mathbf{A}_{t'}) - (\mathbf{B}_t - \mathbf{B}_{t'}) \right\|_F \\ &\leq \|\mathbf{W}\|_F \|\mathbf{A}_t - \mathbf{A}_{t'}\|_F + \|\mathbf{B}_t - \mathbf{B}_{t'}\|_F \\ &\leq \sqrt{K} (\|\mathbf{A}_t - \mathbf{A}_{t'}\|_F + \|\mathbf{B}_t - \mathbf{B}_{t'}\|_F). \end{aligned} \quad (\text{S-36})$$

From (S-36), it is easy to see $\bar{h}(\mathbf{W}) = \mathbf{W}\bar{\mathbf{A}} - \bar{\mathbf{B}}$ since $\sup_{\mathbf{W} \in \mathcal{C}} \|\nabla \tilde{f}_t(\mathbf{W}) - \bar{h}(\mathbf{W})\|_F \leq \sqrt{K} (\|\mathbf{A}_t - \bar{\mathbf{A}}\|_F + \|\mathbf{B}_t - \bar{\mathbf{B}}\|_F)$.

Next, define $g_t \triangleq \tilde{f}_t - f_t$, for all $t \geq 1$. By definition, we have $g_t(\mathbf{W}) \geq 0$, for any $\mathbf{W} \in \mathcal{C}$. Since $\tilde{f}_t \xrightarrow{u} \bar{f}$ and $f_t \xrightarrow{u} f$ (by Glivenko-Cantelli theorem), we have $g_t \xrightarrow{u} \bar{g} \triangleq \bar{f} - f$. Since both \bar{f} and f are differentiable, \bar{g} is differentiable and $\nabla f = \nabla \bar{f} - \nabla \bar{g}$. To show $\bar{\mathbf{W}}$ is a stationary point of f , it suffices to show for any $\mathbf{W} \in \mathcal{C}$, the directional derivative $\langle \nabla f(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle \geq 0$. We show this by proving $\langle \nabla \bar{f}(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle \geq 0$ and $\langle \nabla \bar{g}(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle = 0$ for any $\mathbf{W} \in \mathcal{C}$.

By definition, for any $\mathbf{W} \in \mathcal{C}$ and $t \geq 1$, we have $\tilde{f}_t(\mathbf{W}_t) \leq \tilde{f}_t(\mathbf{W})$. First consider

$$\begin{aligned} \left| \tilde{f}_t(\mathbf{W}_t) - \bar{f}(\bar{\mathbf{W}}) \right| &= \left| \tilde{f}_t(\mathbf{W}_t) - \bar{f}(\mathbf{W}_t) + \bar{f}(\mathbf{W}_t) - \bar{f}(\bar{\mathbf{W}}) \right| \\ &\leq \left\| \tilde{f}_t - \bar{f} \right\|_{\mathcal{C}} + \left| \bar{f}(\mathbf{W}_t) - \bar{f}(\bar{\mathbf{W}}) \right|. \end{aligned}$$

Since $\tilde{f}_t \xrightarrow{u} \bar{f}$ and \bar{f} is continuous, we have $\tilde{f}_t(\mathbf{W}_t) \rightarrow \bar{f}(\bar{\mathbf{W}})$ as $t \rightarrow \infty$. Thus $\bar{f}(\bar{\mathbf{W}}) \leq \bar{f}(\mathbf{W})$, for any $\mathbf{W} \in \mathcal{C}$. This implies $\langle \nabla \bar{f}(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle \geq 0$.

Next we show $\langle \nabla \bar{g}(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle = 0$ for any $\mathbf{W} \in \mathcal{C}$. It suffices to show $\nabla \bar{g}(\bar{\mathbf{W}}) = \mathbf{0}$. Since both \tilde{f}_t and f_t are differentiable, g_t is differentiable and $\nabla g_t = \nabla \tilde{f}_t - \nabla f_t$. First it is easy to see ∇g_t is Lipschitz on \mathcal{C} with constant $L > 0$ independent of t since both $\nabla \tilde{f}_t$ and ∇f_t are Lipschitz on \mathcal{C} with constants independent of t . It is possible to construct another differentiable function \tilde{g}_t with domain $\mathbb{R}^{F \times K}$ such that \tilde{g}_t is nonnegative with a L -Lipschitz gradient on $\mathbb{R}^{F \times K}$ and $\tilde{g}_t(\mathbf{W}) = g_t(\mathbf{W})$ for all $\mathbf{W} \in \mathcal{C}$. Thus

$$\frac{1}{2L} \|\nabla g_t(\mathbf{W}_t)\|_F^2 = \frac{1}{2L} \|\nabla \tilde{g}_t(\mathbf{W}_t)\|_F^2 \leq \tilde{g}_t(\mathbf{W}_t) - \tilde{g}_t^* \leq g_t(\mathbf{W}_t) = \tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t),$$

where $\tilde{g}_t^* = \inf_{\mathbf{W} \in \mathbb{R}^{F \times K}} \tilde{g}_t(\mathbf{W}) \geq 0$. Since $\tilde{f}_t(\mathbf{W}_t) - f_t(\mathbf{W}_t) \rightarrow 0$, we have $\nabla g_t(\mathbf{W}_t) \rightarrow \mathbf{0}$ as $t \rightarrow \infty$. Now consider the first-order Taylor expansion of g_t at \mathbf{W}_t

$$g_t(\mathbf{W}) = g_t(\mathbf{W}_t) + \langle \nabla g_t(\mathbf{W}_t), \mathbf{W} - \mathbf{W}_t \rangle + o(\|\mathbf{W} - \mathbf{W}_t\|_F), \quad \forall \mathbf{W} \in \mathcal{C}. \quad (\text{S-37})$$

As $t \rightarrow \infty$, we have

$$\bar{g}(\mathbf{W}) = \bar{g}(\bar{\mathbf{W}}) + o(\|\mathbf{W} - \bar{\mathbf{W}}\|_F), \quad \forall \mathbf{W} \in \mathcal{C}. \quad (\text{S-38})$$

On the other hand, we have

$$\bar{g}(\mathbf{W}) = \bar{g}(\bar{\mathbf{W}}) + \langle \nabla \bar{g}(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle + o(\|\mathbf{W} - \bar{\mathbf{W}}\|_F), \quad \forall \mathbf{W} \in \mathcal{C}. \quad (\text{S-39})$$

Comparing (S-38) and (S-39), we have

$$\langle \nabla \bar{g}(\bar{\mathbf{W}}), \mathbf{W} - \bar{\mathbf{W}} \rangle + o(\|\mathbf{W} - \bar{\mathbf{W}}\|_F) = 0, \quad \forall \mathbf{W} \in \mathcal{C}. \quad (\text{S-40})$$

Therefore we conclude $\nabla \bar{g}(\bar{\mathbf{W}}) = \mathbf{0}$.

S-7. PROOF OF LEMMA S-1

Since \mathbf{v} , \mathbf{W} , $\mathbf{h}^*(\mathbf{v}, \mathbf{W})$ and $\mathbf{r}^*(\mathbf{v}, \mathbf{W})$ are all bounded, there exist positive constants M_1, M_2, M_3 and M_4 that upper bound $\|\mathbf{v}\|$, $\|\mathbf{W}\|$, $\|\mathbf{h}^*(\mathbf{v}, \mathbf{W})\|$ and $\|\mathbf{r}^*(\mathbf{v}, \mathbf{W})\|$ respectively. Here the matrix norm is the one induced by the (general) vector norm. Take arbitrary \mathbf{W}_1 and \mathbf{W}_2 in \mathcal{C} , we have

$$\begin{aligned} \|\nabla f(\mathbf{W}_1) - \nabla f(\mathbf{W}_2)\| &= \|\mathbb{E}_{\mathbf{v}}[\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}_1) - \nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}_2)]\| \\ &\leq \mathbb{E}_{\mathbf{v}} \|\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}_1) - \nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}_2)\|. \end{aligned}$$

Fix an arbitrary $\mathbf{v} \in \mathcal{V}$ we have

$$\begin{aligned} \|\nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}_1) - \nabla_{\mathbf{W}} \ell(\mathbf{v}, \mathbf{W}_2)\| &\leq \left\| \mathbf{W}_1 \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)^T - \mathbf{W}_2 \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)^T \right\| \\ &\quad + \left\| \mathbf{r}^*(\mathbf{v}, \mathbf{W}_1) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)^T - \mathbf{r}^*(\mathbf{v}, \mathbf{W}_2) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)^T \right\| \\ &\quad + \left\| \mathbf{v} \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)^T - \mathbf{v} \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)^T \right\|. \end{aligned} \quad (\text{S-41})$$

We bound each term on the RHS of (S-41) as follows

$$\begin{aligned}
& \|\mathbf{W}_1 \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)^T - \mathbf{W}_2 \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)^T\| \\
& \leq \|\mathbf{W}_1\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_1) - \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| + \|\mathbf{W}_1 \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1) - \mathbf{W}_2 \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| \\
& \leq M_2 M_3 c_1 \|\mathbf{W}_1 - \mathbf{W}_2\| + (\|\mathbf{W}_1\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_1) - \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| + \|\mathbf{W}_1 - \mathbf{W}_2\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\|) \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| \\
& \leq M_2 M_3 c_1 \|\mathbf{W}_1 - \mathbf{W}_2\| + M_3 (M_2 c_1 \|\mathbf{W}_1 - \mathbf{W}_2\| + M_3 \|\mathbf{W}_1 - \mathbf{W}_2\|) \\
& = (2c_1 M_2 M_3 + M_3^2) \|\mathbf{W}_1 - \mathbf{W}_2\|,
\end{aligned}$$

$$\begin{aligned}
& \|\mathbf{r}^*(\mathbf{v}, \mathbf{W}_1) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)^T - \mathbf{r}^*(\mathbf{v}, \mathbf{W}_2) \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)^T\| \\
& \leq \|\mathbf{r}^*(\mathbf{v}, \mathbf{W}_1)\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_1) - \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| + \|\mathbf{r}^*(\mathbf{v}, \mathbf{W}_1) - \mathbf{r}^*(\mathbf{v}, \mathbf{W}_2)\| \|\mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)\| \\
& \leq c_1 M_4 \|\mathbf{W}_1 - \mathbf{W}_2\| + c_2 M_3 \|\mathbf{W}_1 - \mathbf{W}_2\| \\
& \leq (c_1 M_4 + c_2 M_3) \|\mathbf{W}_1 - \mathbf{W}_2\|,
\end{aligned}$$

and

$$\|\mathbf{v} \mathbf{h}^*(\mathbf{v}, \mathbf{W}_1)^T - \mathbf{v} \mathbf{h}^*(\mathbf{v}, \mathbf{W}_2)^T\| \leq c_1 M_1 \|\mathbf{W}_1 - \mathbf{W}_2\|.$$

Thus, take $c_3 = 2c_1 M_2 M_3 + M_3^2 + c_1 M_4 + c_2 M_3 + c_1 M_1$ and we finish the proof.

S-8. PROOF OF LEMMA S-2

It suffices to show $\|\nabla g(\mathbf{b})\|_2 \leq c_1 \|\mathbf{z} - \mathbf{z}'\|_2 + c_2 \|\mathbf{A} - \mathbf{A}'\|_2$ for any $\mathbf{b} \in \mathcal{B}$ and some positive constants c_1 and c_2 (independent of \mathbf{b}). We write $\|\nabla g(\mathbf{b})\|_2$ as

$$\begin{aligned}
\|\nabla g(\mathbf{b})\|_2 &= \|\mathbf{A}'^T (\mathbf{A}' \mathbf{b} - \mathbf{z}') - \mathbf{A}^T (\mathbf{A} \mathbf{b} - \mathbf{z})\|_2 \\
&= \|(\mathbf{A}'^T \mathbf{A}' - \mathbf{A}^T \mathbf{A}) \mathbf{b} - (\mathbf{A}'^T \mathbf{z}' - \mathbf{A}^T \mathbf{z})\|_2 \\
&\leq \underbrace{\|\mathbf{A}'^T \mathbf{A}' - \mathbf{A}^T \mathbf{A}\|_2}_{\langle 1 \rangle} \|\mathbf{b}\|_2 + \underbrace{\|\mathbf{A}'^T \mathbf{z}' - \mathbf{A}^T \mathbf{z}\|_2}_{\langle 2 \rangle}.
\end{aligned}$$

By the compactness of \mathcal{Z} , \mathcal{A} and \mathcal{B} , there exist positive constants M_1 , M_2 and M_3 such that $\|\mathbf{z}\|_2 \leq M_1$, $\|\mathbf{A}\|_2 \leq M_2$ and $\|\mathbf{b}\|_2 \leq M_3$, for any $\mathbf{z} \in \mathcal{Z}$, $\mathbf{A} \in \mathcal{A}$ and $\mathbf{b} \in \mathcal{B}$. Thus,

$$\begin{aligned}
\langle 1 \rangle &\leq M_3 \|\mathbf{A}^T \mathbf{A} - \mathbf{A}^T \mathbf{A}' + \mathbf{A}^T \mathbf{A}' - \mathbf{A}'^T \mathbf{A}'\|_2 \\
&\leq M_3 (\|\mathbf{A}^T (\mathbf{A} - \mathbf{A}')\|_2 + \|(\mathbf{A} - \mathbf{A}')^T \mathbf{A}'\|_2) \\
&\leq 2M_2 M_3 \|\mathbf{A} - \mathbf{A}'\|_2.
\end{aligned}$$

Similarly for $\langle 2 \rangle$ we have

$$\begin{aligned}
\langle 2 \rangle &= \|\mathbf{A}^T \mathbf{z} - \mathbf{A}^T \mathbf{z}' + \mathbf{A}^T \mathbf{z}' - \mathbf{A}'^T \mathbf{z}'\|_2 \\
&\leq \|\mathbf{A}^T (\mathbf{z} - \mathbf{z}')\|_2 + \|(\mathbf{A} - \mathbf{A}')^T \mathbf{z}'\|_2 \\
&\leq M_2 \|\mathbf{z} - \mathbf{z}'\|_2 + M_1 \|\mathbf{A} - \mathbf{A}'\|_2.
\end{aligned}$$

Hence $\langle 1 \rangle + \langle 2 \rangle \leq M_2 \|\mathbf{z} - \mathbf{z}'\|_2 + (M_1 + 2M_2 M_3) \|\mathbf{A} - \mathbf{A}'\|_2$. We now take $c_1 = M_2$ and $c_2 = M_1 + 2M_2 M_3$ to complete the proof.

S-9. TECHNICAL LEMMAS

Lemma S-3 ([65, Lemma 5]). *Let I be a closed interval in \mathbb{R} . Define $g_{\tau, I}(t) = \tau |t| + \delta_I(t)$, where δ_I is the indicator function for the interval I . Then the proximal operator for $g_{\tau, I}$ is given by*

$$\text{prox}_{g_{\tau, I}}(q) = \Pi_I(\mathcal{S}_\tau(q)), \quad (\text{S-42})$$

where $q \in \mathbb{R}$, \mathcal{S}_τ is the soft-thresholding operator with threshold τ and Π_I is the Euclidean projector onto the interval I .

Lemma S-4 (Projection onto nonnegative ℓ_2 balls). *Let $\mathcal{C}' \triangleq \{\mathbf{x} \in \mathbb{R}_+^n \mid \|\mathbf{x}\| \leq 1\}$. Then for all $\mathbf{y} \in \mathbb{R}^n$,*

$$\Pi_{\mathcal{C}'}(\mathbf{y}) = \frac{\mathbf{y}_+}{\max\{1, \|\mathbf{y}_+\|_2\}}, \quad (\text{S-43})$$

where $(\mathbf{y}_+)_i = \max\{0, y_i\}$, $\forall i \in [n]$.

Proof: The KKT conditions for

$$\begin{aligned} \min \quad & \|\mathbf{y} - \mathbf{x}\|_2^2 \\ \text{s. t.} \quad & \mathbf{x} \geq 0, \|\mathbf{x}\|_2 \leq 1 \end{aligned}$$

are given by

$$\mathbf{x}^* \geq 0, \|\mathbf{x}^*\|_2 \leq 1, \lambda^* \geq 0 \quad (\text{S-44})$$

$$(\lambda^* + 1)\mathbf{x}^* \geq \mathbf{y}, \quad (\text{S-45})$$

$$\lambda^*(\|\mathbf{x}^*\|_2^2 - 1) = 0, \quad (\text{S-46})$$

$$(\lambda^* + 1)(x_i^*)^2 = x_i^* y_i, \quad \forall i \in [n]. \quad (\text{S-47})$$

We fix an $i \in [n]$. Define $\mathcal{I} \triangleq \{i \in [n] \mid y_i > 0\}$. For any $\mathbf{z} \in \mathbb{R}^n$, let $\mathbf{z}_{\mathcal{I}}$ be the subvector of \mathbf{z} with indices from \mathcal{I} . If $y_i \leq 0$, then by (S-47) $x_i^* = 0$. If $y_i > 0$, then by (S-45) $x_i^* > 0$. Thus by (S-47) we have $x_i^* = y_i/(\lambda^* + 1)$. If $\lambda^* = 0$, then $x_i^* = y_i$. In this case $\|\mathbf{y}_{\mathcal{I}}\|_2 = \|\mathbf{x}_{\mathcal{I}}^*\|_2 = \|\mathbf{x}^*\|_2 \leq 1$. If $\lambda^* > 0$, then by (S-46) $\|\mathbf{x}^*\|_2^2 = 1$. Then $\|\mathbf{x}^*\|_2^2 = \|\mathbf{x}_{\mathcal{I}}^*\|_2^2 = \|\mathbf{y}_{\mathcal{I}}\|_2^2/(\lambda^* + 1)^2 = 1$. This means $\lambda^* + 1 = \|\mathbf{y}_{\mathcal{I}}\|_2$ so $x_i^* = y_i/\|\mathbf{y}_{\mathcal{I}}\|_2$. Also, we notice in such case $y_i > x_i^* > 0$ so $\|\mathbf{y}_{\mathcal{I}}\|_2 > 1$. Combining both cases where $\lambda^* = 0$ and $\lambda^* > 0$, we have $x_i^* = y_i/\max\{1, \|\mathbf{y}_{\mathcal{I}}\|_2\}$, for all $i \in \mathcal{I}$. ■

Lemma S-5 (Danskin's Theorem; [98, Theorem 4.1]). *Let \mathcal{X} be a metric space and \mathcal{U} be a normed vector space. Let $f : \mathcal{X} \times \mathcal{U} \rightarrow \mathbb{R}$ have the following properties*

- 1) $f(x, \cdot)$ is differentiable on \mathcal{U} , for any $x \in \mathcal{X}$.
- 2) $f(x, u)$ and $\nabla_u f(x, u)$ are continuous on $\mathcal{X} \times \mathcal{U}$.

Let Φ be a compact set in \mathcal{X} . Define $v(u) = \inf_{x \in \Phi} f(x, u)$ and $S(u) = \arg \min_{x \in \Phi} f(x, u)$, then $v(u)$ is directionally differentiable and its directional derivative along $d \in \mathcal{U}$, $v'(u, d)$ is given by

$$v'(u, d) = \min_{x \in S(u)} \nabla_u f(x, u)^T d. \quad (\text{S-48})$$

In particular, if for some $u_0 \in \mathcal{U}$, $S(u_0) = \{x_0\}$, then v is differentiable at $u = u_0$ and $\nabla v(u_0) = \nabla_u f(x_0, u_0)$.

Lemma S-6 (The Maximum Theorem; [99, Theorem 14.2.1 & Example 2]). *Let \mathcal{P} and \mathcal{X} be two metric spaces. Consider a maximization problem*

$$\max_{x \in B(p)} f(p, x), \quad (\text{S-49})$$

where $B : \mathcal{P} \rightarrow \mathcal{X}$ is a correspondence and $f : \mathcal{P} \times \mathcal{X} \rightarrow \mathbb{R}$ is a function. If B is compact-valued and continuous on \mathcal{P} and f is continuous on $\mathcal{P} \times \mathcal{X}$, then the correspondence $S(p) = \arg \max_{x \in B(p)} f(p, x)$ is compact-valued and upper hemicontinuous, for any $p \in \mathcal{P}$. In particular, if for some $p_0 \in \mathcal{P}$, $S(p_0) = \{s(p_0)\}$, where $s : \mathcal{P} \rightarrow \mathcal{X}$ is a function, then s is continuous at $p = p_0$. Moreover, we have the same conclusions if the maximization in (S-49) is replaced by minimization.

Lemma S-7 (Leibniz Integral Rule). *Let \mathcal{X} be an open set in \mathbb{R}^n and let $(\Omega, \mathcal{A}, \mu)$ be a measure space. If $f : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$ satisfies*

- 1) For all $x \in \mathcal{X}$, the mapping $\omega \mapsto f(x, \omega)$ is Lebesgue integrable.
- 2) For all $\omega \in \Omega$, $\nabla_x f(x, \omega)$ exists on \mathcal{X} .
- 3) For all $x \in \mathcal{X}$, the mapping $\omega \mapsto \nabla_x f(x, \omega)$ is Lebesgue integrable.

Then $\int_{\Omega} f(x, \omega) d\mu(\omega)$ is differentiable and

$$\nabla_x \int_{\Omega} f(x, \omega) d\mu(\omega) = \int_{\Omega} \nabla_x f(x, \omega) d\mu(\omega). \quad (\text{S-50})$$

Remark 7. This is a simplified version of the Leibniz Integral Rule. See [100, Theorem 16.8] for weaker conditions on f .

Definition S-1 (Quasi-martingale; [101, Definition 1.4]). Let $\{X_t\}_{t \in \mathcal{T}}$ be a stochastic process and let $\{\mathcal{F}_t\}_{t \in \mathcal{T}}$ be the filtration to which $\{X_t\}_{t \in \mathcal{T}}$ is adapted, where \mathcal{T} is a subset of the real line. We call $\{X_t\}_{t \in \mathcal{T}}$ a quasi-martingale if there exist two stochastic processes $\{Y_t\}_{t \in \mathcal{T}}$ and $\{Z_t\}_{t \in \mathcal{T}}$ such that

- i) both $\{Y_t\}_{t \in \mathcal{T}}$ and $\{Z_t\}_{t \in \mathcal{T}}$ are adapted to $\{\mathcal{F}_t\}_{t \in \mathcal{T}}$,
- ii) $\{Y_t\}_{t \in \mathcal{T}}$ is a martingale and $\{Z_t\}_{t \in \mathcal{T}}$ has bounded variations on \mathcal{T} a.s.,
- iii) $X_t = Y_t + Z_t$, for all $t \in \mathcal{T}$ a.s..

Lemma S-8 (The Quasi-martingale Convergence Theorem; [41, Theorem 9.4 & Proposition 9.5]). *Let $(u_t)_{t \geq 1}$ be a nonnegative discrete-time stochastic process on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, i.e., $u_t \geq 0$ a.s., for all $t \geq 1$. Let $\{\mathcal{F}_t\}_{t \geq 1}$ be a filtration*

to which $(u_t)_{t \geq 1}$ is adapted. Define another binary stochastic process $(\delta_t)_{t \geq 1}$ as

$$\delta_t = \begin{cases} 1, & \text{if } \mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t] > 0 \\ 0, & \text{otherwise} \end{cases}. \quad (\text{S-51})$$

If $\sum_{t=1}^{\infty} \mathbb{E}[\delta_t(u_{t+1} - u_t)] < \infty$, then $u_t \xrightarrow{\text{a.s.}} u$, where u is integrable on $(\Omega, \mathcal{F}, \mathbb{P})$ and nonnegative a.s.. Furthermore, $(u_t)_{t \geq 1}$ is a quasi-martingale and

$$\sum_{t=1}^{\infty} \mathbb{E}[|\mathbb{E}[u_{t+1} - u_t | \mathcal{F}_t]|] < \infty. \quad (\text{S-52})$$

Lemma S-9 (Donsker's Theorem; [40, Section 19.2]). Let X_1, \dots, X_n be i.i.d. generated from a distribution \mathbb{P} . Define the empirical distribution $\mathbb{P}_n \triangleq \frac{1}{n} \sum_{i=1}^n \delta_{X_i}$. For a measurable function f , define $\mathbb{P}_n f$ and $\mathbb{P} f$ as the expectations of f under the distributions \mathbb{P}_n and \mathbb{P} respectively. Define an empirical process $G_n(f) \triangleq \sqrt{n}(\mathbb{P}_n f - \mathbb{P} f)$, $f \in \mathcal{F}$, where \mathcal{F} is a class of measurable functions. \mathcal{F} is \mathbb{P} -Donsker if and only if the sequence of empirical processes $\{G_n\}_{n \geq 1}$ converges in distribution to a zero-mean Gaussian process G tight in $\ell^\infty(\mathcal{F})$, where $\ell^\infty(\mathcal{F})$ is the space of all real-valued and bounded functionals defined on \mathcal{F} equipped with the uniform norm on \mathcal{F} , denoted as $\|\cdot\|_{\mathcal{F}}$. Moreover, in such case, we have $\mathbb{E} \|G_n\|_{\mathcal{F}} \rightarrow \mathbb{E} \|G\|_{\mathcal{F}}$.

Lemma S-10 (Glivenko-Cantelli theorem; [40, Section 19.2]). Let \mathbb{P}_n and \mathbb{P} be the distributions defined as in Lemma S-9. A class of measurable functions \mathcal{F} is \mathbb{P} -Glivenko-Cantelli if and only if $\sup_{f \in \mathcal{F}} |\mathbb{P}_n f - \mathbb{P} f| \xrightarrow{\text{a.s.}} 0$.

Lemma S-11 (A sufficient condition for \mathbb{P} -Glivenko-Cantelli and \mathbb{P} -Donsker classes; [40, Example 19.7]). Define a probability space $(\mathcal{X}, \mathcal{A}, \mathbb{P})$. Let $\mathcal{F} = \{f_\theta : \mathcal{X} \rightarrow \mathbb{R} \mid \theta \in \Theta\}$ be a class of measurable functions, where Θ is a bounded subset in \mathbb{R}^d . If there exists a universal constant $K > 0$ such that

$$|f_{\theta_1}(x) - f_{\theta_2}(x)| \leq K \|\theta_1 - \theta_2\|, \quad \forall \theta_1, \theta_2 \in \Theta, \quad \forall x \in \mathcal{X}, \quad (\text{S-53})$$

where $\|\cdot\|$ is a general vector norm in \mathbb{R}^d , then \mathcal{F} is both \mathbb{P} -Glivenko-Cantelli and \mathbb{P} -Donsker.

Lemma S-12 ([44, Lemma 8]). Let $(a_n), (b_n)$ be two nonnegative sequences. Suppose $\sum_{n=1}^{\infty} a_n = \infty$ and $\sum_{n=1}^{\infty} a_n b_n < \infty$, and $\exists N \in \mathbb{N}$ and $K > 0$ such that for all $n \geq N$, $|b_{n+1} - b_n| \leq K a_n$. Then (b_n) converges and $\lim_{n \rightarrow \infty} b_n = 0$.

S-10. ADDITIONAL EXPERIMENT RESULTS

This section consists of two parts. In the first part, we show the convergence speeds of all the online and batch algorithms on the (contaminated) CBCL face dataset for different values of the mini-batch size τ , the latent dimension K , the penalty parameter ρ in the ADMM-based algorithms, the step-size parameter κ in the PGD-based algorithms and the (salt and pepper) noise density parameters ν and $\tilde{\nu}$ in Figure 10 to 15. As mentioned in Section VII-D, all the convergence results on the CBCL face dataset agree with those on the synthetic dataset. In the second part, we show the quality of the denoised images of all the algorithms on the CBCL face dataset for $K = 25$ and $K = 100$ in Table IV (a) and (b) respectively. The corresponding running times of all the algorithms for $K = 25$ and $K = 100$ are shown in Table V (a) and (b) respectively. As stated in Section VII-F, the results when $K = 25$ and $K = 100$ are similar to those when $K = 49$.

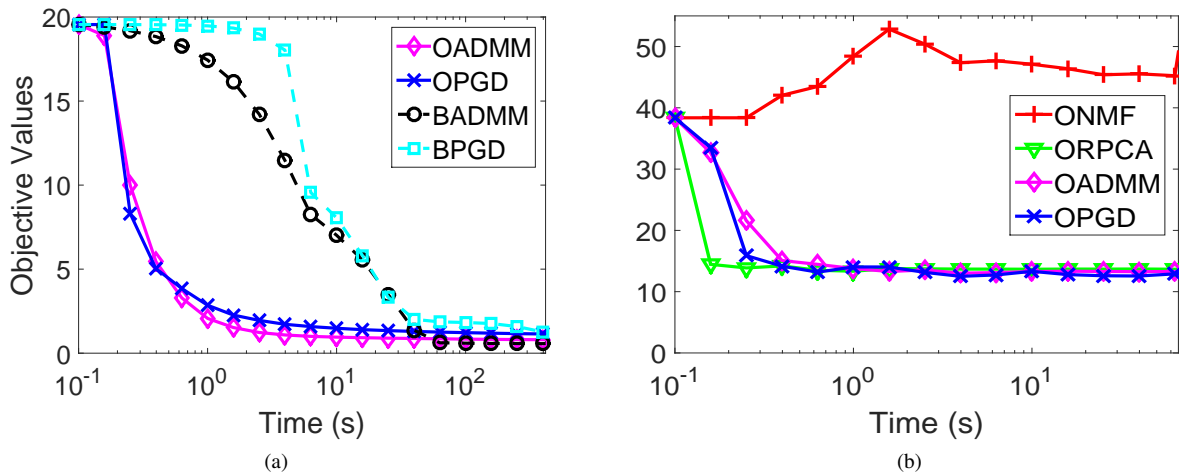


Fig. 10. The objective values (as a function of time) of (a) our online algorithms and their batch counterparts (b) our online algorithms and other online algorithms on the CBCL face dataset. The parameters are set according to the canonical setting.

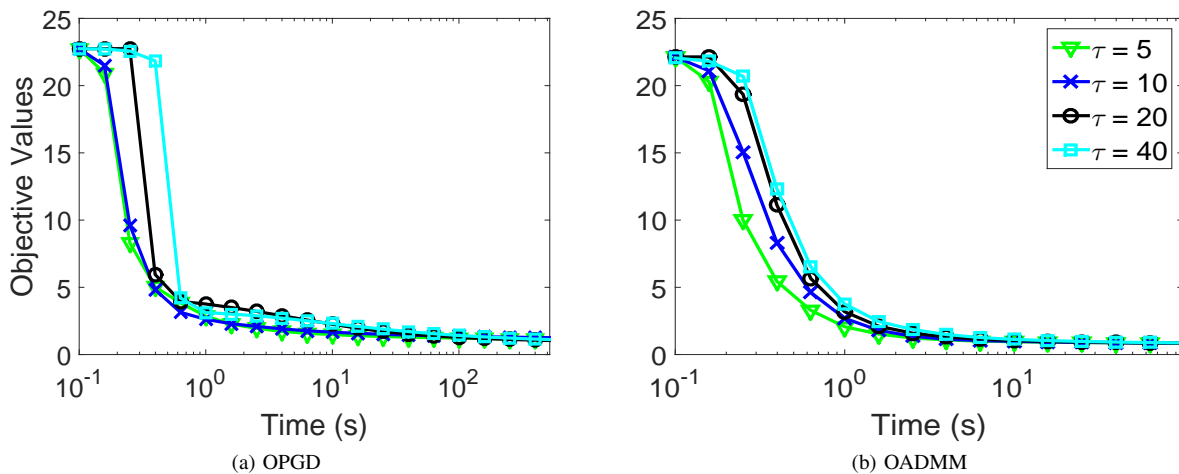


Fig. 11. The objective values (as a function of time) of (a) OPGD and (b) OADMM for different values of τ on the CBCL face dataset. All the other parameters are set according to the canonical setting.

TABLE IV
PSNRs (IN DB) OF ALL THE ALGORITHMS ON THE CBCL FACE DATASET WITH DIFFERENT NOISE DENSITY.

	Setting 1	Setting 2	Setting 3
OADMM	11.39 ± 0.16	11.37 ± 0.12	11.35 ± 0.18
OPGD	11.49 ± 0.05	11.43 ± 0.09	11.38 ± 0.06
BADMM	11.51 ± 0.19	11.46 ± 0.07	11.41 ± 0.15
BPGD	11.54 ± 0.07	11.46 ± 0.17	11.42 ± 0.15
ONMF	5.99 ± 0.04	5.97 ± 0.12	5.97 ± 0.08
ORPCA	11.26 ± 0.05	11.24 ± 0.11	11.22 ± 0.11

(a) $K = 25$

	Setting 1	Setting 2	Setting 3
OADMM	11.39 ± 0.02	11.35 ± 0.11	11.35 ± 0.06
OPGD	11.51 ± 0.01	11.45 ± 0.09	11.44 ± 0.11
BADMM	11.51 ± 0.11	11.47 ± 0.00	11.45 ± 0.03
BPGD	11.52 ± 0.16	11.46 ± 0.07	11.45 ± 0.12
ONMF	5.99 ± 0.19	5.97 ± 0.03	5.95 ± 0.05
ORPCA	11.26 ± 0.03	11.24 ± 0.16	11.20 ± 0.13

(b) $K = 100$

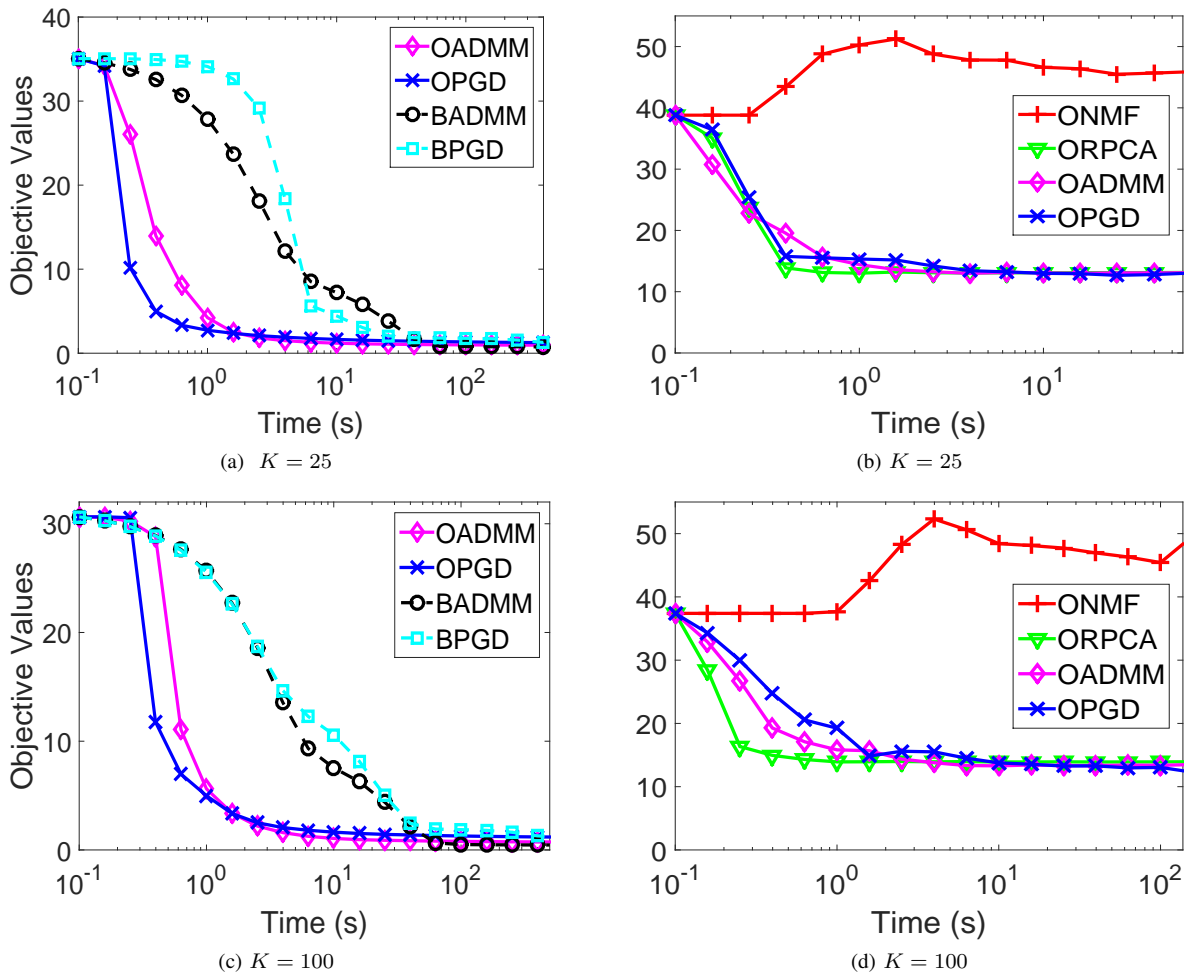


Fig. 12. The objective values (as a function of time) of all the algorithms for different values of K on the CBCL face dataset. In (a) and (b), $K = 25$. In (c) and (d), $K = 100$. All the other parameters are set according to the canonical setting.

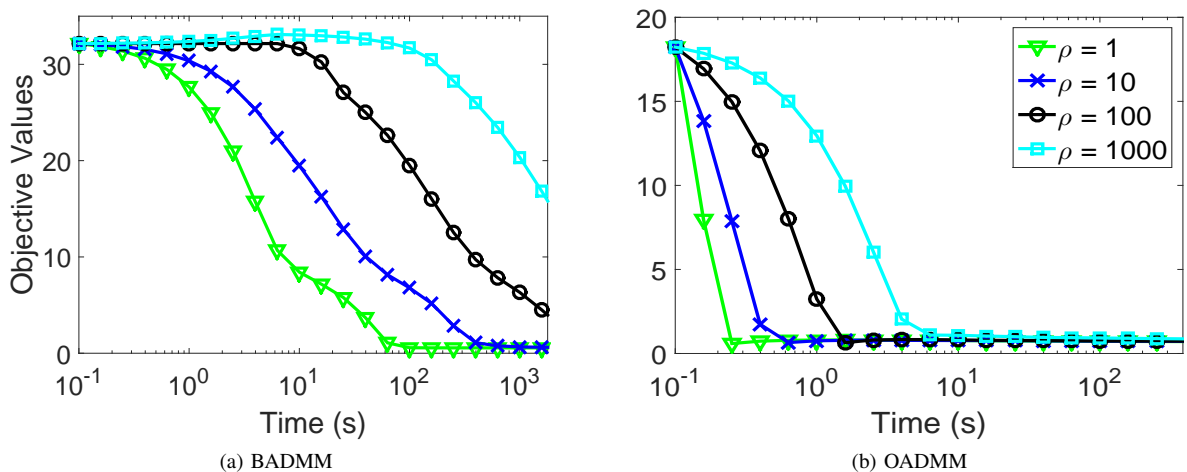


Fig. 13. The objective values (as a function of time) of (a) BADMM and (b) OADMM for different values of ρ on the CBCL face dataset. All the other parameters are set according to the canonical setting.

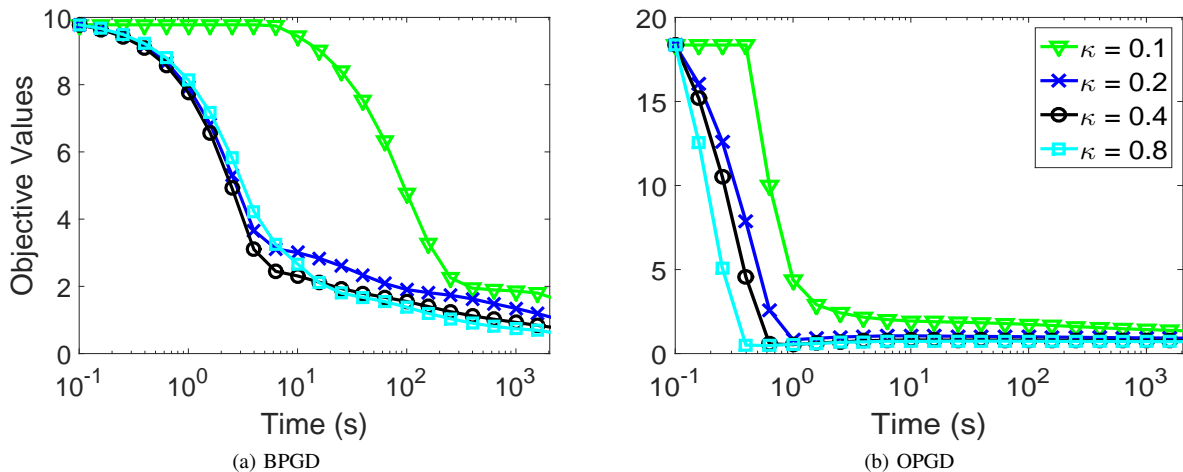


Fig. 14. The objective values (as a function of time) of (a) BPGD and (b) OPGD for different values of κ on the CBCL face dataset. All the other parameters are set according to the canonical setting.

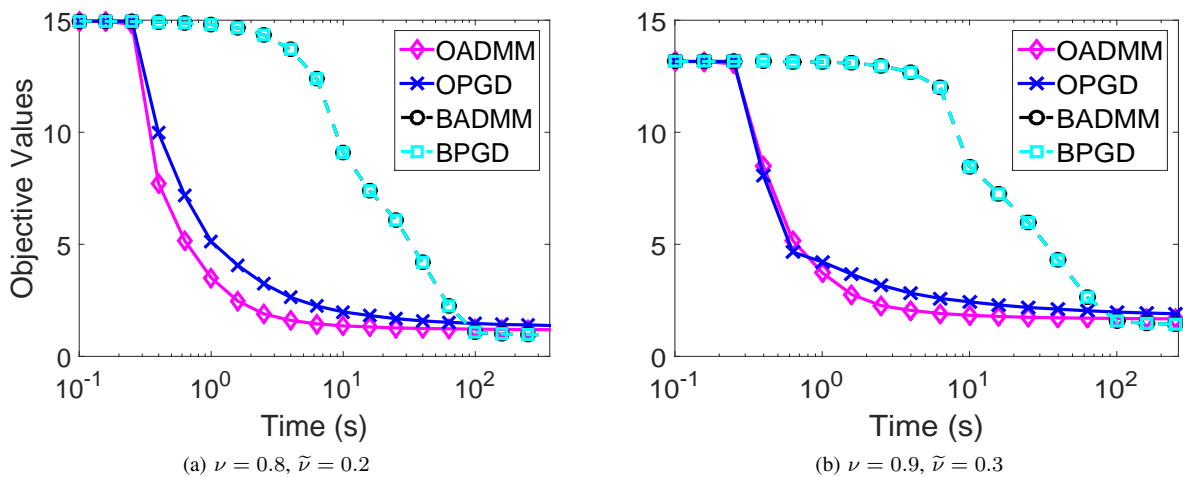


Fig. 15. The objective values (as a function of time) of our online and batch algorithms on the CBCL face dataset with a larger proportion of outliers.

TABLE V
RUNNING TIMES (IN SECONDS) OF ALL THE ALGORITHMS ON THE CBCL FACE DATASET WITH DIFFERENT NOISE DENSITY.

	Setting 1	Setting 2	Setting 3
OADMM	420.58 \pm 2.59	427.66 \pm 4.80	430.02 \pm 2.93
OPGD	431.66 \pm 2.49	455.15 \pm 1.70	463.67 \pm 1.12
BADMM	1009.45 \pm 11.27	1184.29 \pm 10.49	1240.91 \pm 8.21
BPGD	1125.58 \pm 12.83	1185.64 \pm 13.36	1279.07 \pm 9.08
ONMF	2384.70 \pm 9.59	2588.29 \pm 14.39	2698.57 \pm 10.24
ORPCA	365.98 \pm 5.29	382.49 \pm 4.20	393.10 \pm 4.27

(a) $K = 25$

	Setting 1	Setting 2	Setting 3
OADMM	422.97 \pm 2.38	424.67 \pm 2.65	434.17 \pm 4.67
OPGD	430.19 \pm 2.27	448.72 \pm 3.90	454.30 \pm 4.65
BADMM	1009.04 \pm 8.53	1187.58 \pm 5.06	1250.53 \pm 4.67
BPGD	1131.89 \pm 7.04	1192.46 \pm 7.43	1280.55 \pm 7.93
ONMF	2379.86 \pm 15.18	2591.09 \pm 11.91	2693.08 \pm 12.48
ORPCA	363.37 \pm 3.01	390.58 \pm 5.31	401.21 \pm 3.27

(b) $K = 100$