# Low-complexity Proximal Gauss-Newton Algorithm for Nonnegative Matrix Factorization

Kejun Huang Department of CISE University of Florida kejun.huang@ufl.edu

Abstract-In this paper we propose a quasi-Newton algorithm for the celebrated nonnegative matrix factorization (NMF) problem. The proposed algorithm falls into the general framework of Gauss-Newton and Levenberg-Marquardt methods. However, these methods were not able to handle constraints, which is present in NMF. One of the key contributions in this paper is to apply alternating direction method of multipliers (ADMM) to obtain the iterative update from this Gauss-Newton-like algorithm. Furthermore, we carefully study the structure of the Jacobian Gramian matrix given by the Gauss-Newton updates, and designed a way of exactly inverting the matrix with complexity O(mnk), which is a significant reduction compared to the naive implementation of complexity  $O((m + n)^3 k^3)$ . The resulting algorithm, which we call NLS-ADMM, enjoys fast convergence rate brought by the quasi-Newton algorithmic framework, while maintaining low per-iteration complexity similar to that of alternating algorithms. Numerical experiments on synthetic data confirms the efficiency of our proposed algorithm.

*Index Terms*—nonnegative matrix factorization, Gauss-Newton, proximal algorithm, alternating direction method of multipliers, lowcomplexity

## I. INTRODUCTION

Nonnegative matrix factorization, or NMF for short, is a powerful data analysis tool for signal processing and machine learning [1]. First popularized by Lee and Seung [2] by showing its capability of "learning the parts of objects" in computer vision, NMF has found numerous applications by providing meaningful and interpretable results, including blind source separation [3], [4], hyperspectral unmixing [5], [6], topic modeling [7], [8], [9], and community detection [10], [11], to name just a few. In terms of model identifiability, NMF has been shown to be essentially unique under the "sufficiently scattered" condition [12], [13], which explained the wide applicability of NMF.

Computationally, NMF is shown to be NP-hard [14]. Most existing NMF algorithms start by formulating an optimization problem using the least-squares loss between the data matrix  $X \in \mathbb{R}^{m \times n}$  and the product of the two factor matrices  $W \in \mathbb{R}^{m \times k}$  and  $H \in \mathbb{R}^{n \times k}$  with a target rank *k* as follows

$$\min_{\boldsymbol{W} \ge 0, \boldsymbol{H} \ge 0} \|\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H}^{\mathsf{T}}\|_{F}^{2}.$$
(1)

Since (1) is non-convex, there is no algorithm that guarantees to optimally solve it in polynomial-time. The prevailing method is to update the two factors W and H in an alternating fashion. Fixing one of the factors, problem (1) with respect to the other factor becomes a nonnegative least-squares problem, which is convex, but does not admit a closed-form solution. Some well-known examples include:

• Multiplicative update (MU). Lee and Seung [15] proposed to iteratively minimize an auxilirary function that majorizes the

X. Fu is supported in part by the National Science Foundation under Project NSF ECCS-1608961, ECCS 1808159, III-1910118, and the Army Research Office (ARO) under Proejct ARO W911NF-19-1-0247.

Xiao Fu School of EECS Oregon State University xiao.fu@oregonstate.edu

nonnegative least-squares loss, which results in a low-complexity closed-form update;

- Hierarchical alternating least squares (HALS) updates the factors column by column [16];
- Alternating projected gradient (APG) goes one projected gradient step for each of the factor updates [17]; an accelerated version was proposed by taking extrapolations [18];
- Alternating nonnegative least squares with active set method (ANLS-AS) or block principal pivoting (ANLS-BPP). To solve the nonnegative least-squares sub-problems exactly, classical methods like the active-set algorithm and block principal pivoting were applied in an alternating fashion [19], [20];
- Alternating optimization with alternating direction method of multipliers (AO-ADMM). It turns out that nonnegative leastsquares can be solved with essentially the same computational complexity as its unconstrained version using ADMM, even though it does not admit a closed-form solution; the resulting AO-ADMM algorithm is not only efficient but also flexible in incorporating other kinds of constraints/regularizations seamlessly [21].

For the plain NMF problem of minimizing (1), the literature shows that ANLS-BPP [20], APG with extrapolation [18], and AO-ADMM [21] give the best emprical performance.

A few attempts have been made on designing "all-at-once" algorithms to update the two factors simultaneously. The simplest one is perhaps the projected gradient algorithm for the entire problem, which was investigated in the original APG paper [17] and reported to work inferior than the alternating version. ADMM was also applied to the entire problem rather than the convex sub-problems [22], [23]; however, ADMM for non-convex problems is not even guaranteed to monotonically decrease the loss, and in practice do perform less stably than alternating updates.

In this paper, we propose to explore a Gauss-Newton-like algorithm for NMF, which also falls into the "all-at-once" update-rule category. To the best of our knowledge, this classical quasi-Newton algorithm has not been applied to NMF due to two reasons: on the one hand, there are in total (m + n)k variables, and a second-order algorithm without exploiting problem structures would result in a periteration complexity of  $O((m + n)^3k^3)$ , which is prohibitive for most practical problem sizes; on the other hand, the original Gauss-Newton algorithm was designed for *unconstrained* nonlinear least-squares, with the premise that unconstrained linear least-squares problems are easy to solve, whereas in NMF there are nonnnegativity constraints.

**Contributions.** We propose a Gauss-Newton-like algorithm called NLS-ADMM for NMF, which resolves the two aforementioned challenges that hinder the use of such algorithms for NMF. We start by writing out explicitly the highly structured Hessian-approximation

given by the Gauss-Newton method, and show that it can be inverted with complexity O(mnk) by exploiting the structure. Notice that the complexity is significantly lower than directly inverting with complexity  $O((m+n)^3k^3)$ , and is in fact in the same order as a single factor update. There is also no additional memory overhead in storing the Hessian-approximation, thanks again to the nice structures. The nonnegativity constraints are handled by the use of ADMM, so that all the problem structures exploited in solving the unconstrained problem can be seamlessly applied here, with the flexibility of incorporating not only nonnegativity but also a variety of constraints/regularizations with efficient projection/proximity operators.

**Related works.** The closest work is the use of Gauss-Newton and its proximal version, Levenberg-Marquardt, to the canonical polyadic decomposition (CPD) for tensors [24], [25]. It has been shown that the Hessian approximation exhibits similar structures to that of the matrix case [26]. However, because of the subtle differences, people have not been able to bring down the per-iteration complexity down to that of a single factor update like this paper. Furthermore, their ways of handling nonnegativity constraints are either squaring the variables or addition log-barriers to the objective function. However, both approaches destroy the structure of the Hessian approximation for efficient inversion, unless an indirect method like the conjugate gradient method is used [24]. In addition, these approaches are not amendable for other kinds of constraints/regularizations.

#### II. PROPOSED ALGORITHM

#### A. General framework: NLS-ADMM

Before we introduce the general algorithmic framework for our proposed algorithm, we briefly review the general idea of Gauss-Newton and Levenberg-Marquardt algorithm. Details can be found in many textbooks, e.g. [27].

Consider the (unconstrained) nonlinear least-squares problem

minimize 
$$\|f(z)\|^2$$
, (2)

where f(z) is a vector of nonlinear functions with respect to the variable z.

The Gauss-Newton algorithm is an iterative heuristic for solving (2). At iteration t when the update is at  $z_t$ , we first take a linear approximation of f(z) at  $z_t$  as

$$f(z) \approx f(z_t) + J_f(z_t)(z - z_t),$$

where  $J_f(z_t)$  is the Jacobian matrix of f at  $z_t$ , which we will simplify its notation as  $J_t$  in the sequel, with its (p, q)-th entry defined as

$$[\boldsymbol{J}_t]_{p,q} = \left[\boldsymbol{J}_f(z_t)\right]_{p,q} = \frac{\partial f_p}{\partial z_q}(z_t).$$

The update rule of Gauss-Newton is therefore

$$z_{t+1} = \arg\min_{z} \|f(z_t) + J_t(z - z_t)\|^2.$$
(3)

If  $J_t$  has full colum-rank, then the update (3) has a closed form  $x_{t+1} = z_t - (J_t^{\mathsf{T}} J_t)^{-1} J_t^{\mathsf{T}} f(z_t)$ . In the vicinity of the solution, Gauss-Newton shows super-linear convergence rate.

There are some short-comings of Gauss-Newton. It is not a descent algorithm, so if the initialization is not close to the optimal solution, it may take some iterations to even start decreasing the loss function. Furthermore, if  $J_t$  does not have full column-rank, which indeed often happens in practice, the update rull (3) is not even well-defined.

A remedy is to add a proximal term to (3), leading to the Levenberg-Marquardt algorithm

$$z_{t+1} = \arg\min_{z} \|f(z_t) + J_t(z - z_t)\|^2 + \lambda_t \|z - z_t\|^2.$$
(4)

With the help of the proximal term, the update rule is always well defined  $\mathbf{x}_{t+1} = \mathbf{z}_t - (\mathbf{J}_t^{\mathsf{T}} \mathbf{J}_t + \lambda_t \mathbf{I})^{-1} \mathbf{J}_t^{\mathsf{T}} f(z_t)$ . Furthermore, one can adaptively change the value of  $\lambda_t$  to guarantee monotonic decrease of the loss function. One effective approach, which we adopt in this paper, is as follows: if the new update does not decrease the loss, then double  $\lambda_t$  and re-calculate the update; otherwise, accept the update and let  $\lambda_{t+1} = \lambda_t/2$ .

To incorporate constraints to the variable z, say  $z \in \mathbb{Z}$ , we propose to keep the general form of Levenberg-Marquardt update as

$$z_{t+1} = \arg\min_{z \in \mathcal{Z}} ||f(z_t) + J_t(z - z_t)||^2 + \lambda_t ||z - z_t||^2.$$
(5)

We know that the constrained least-squares problem does not admit a closed-form solution in general. However, it can be efficiently calculated via the following ADMM iterates:

$$\widetilde{z} \leftarrow z_t - \left( \boldsymbol{J}_t^{\mathsf{T}} \boldsymbol{J}_t + (\rho + \lambda_t) \boldsymbol{I} \right)^{-1} \left( \boldsymbol{J}_t^{\mathsf{T}} \boldsymbol{f}(z_t) + \rho(z_t - z + \boldsymbol{u}) \right)$$

$$z \leftarrow \operatorname{Proj}_{\mathcal{Z}} (\widetilde{z} + \boldsymbol{u}) \tag{6}$$

$$\boldsymbol{u} \leftarrow \boldsymbol{u} + \widetilde{z} - z$$

where  $\tilde{z}$  is an auxiliary variable subject to the additional constraint  $\tilde{z} = z$ , u is the scaled version of its Lagrange multiplier,  $\rho$  is some positive scalar, and we assume that the projection onto the set  $\mathcal{Z}$  is easy to evaluate. For convex problems, ADMM always converges to a global solution, although the convergence rate may depend on the choice of  $\rho$ . Details on ADMM can be found in [28]. Once we find an accurate enough update of  $z_{t+1}$  per (5) via ADMM, all the other steps remain the same as the classical Levenberg-Marquardt algorithm.

## B. Efficient inversion of the Jacobian Gramian

In the context of NMF, we rearrange the variables and nonlinear functions as follows

$$z = \begin{bmatrix} \operatorname{vec}(W) \\ \operatorname{vec}(H) \end{bmatrix}, \qquad f(z) = \operatorname{vec}(WH^{\top} - X).$$

As we can see from the key computational steps in (6), the computationally heaviest step is to calculate the inverse of the Jacobian Gramian. Since  $z \in \mathbb{R}^{(m+n)k}$ , a naive implementation would require  $\mathcal{O}((m+n)^3k^3)$  flops to compute the inverse, which is prohibitive in most practical cases.

In the supplementary material of [26], the Cramer-Rao bound for matrix factorization was derived. It turns out that replacing the ground-truth W and H with the current update  $W_t$  and  $H_t$  is exactly the Jacobian Gramian, which is highly structured. To further simplify notations, we drop the subscript t when deriving the inverses, and we have [26]

$$J = \begin{bmatrix} H \otimes I_m & C_{nm}(W \otimes I_n) \end{bmatrix},$$
  
$$J^{\mathsf{T}}J = \begin{bmatrix} H^{\mathsf{T}}H \otimes I_m & (I_k \otimes W)C_{kk}(I_k \otimes H)^{\mathsf{T}} \\ (I_k \otimes H)C_{kk}(I_k \otimes W)^{\mathsf{T}} & W^{\mathsf{T}}W \otimes I_n \end{bmatrix}$$

where  $\otimes$  denotes the Kronecker product and  $C_{nm}$  is a special permutation matrix called the commutation matrix, which has the property that  $C_{nm} \operatorname{vec}(S) = \operatorname{vec}(S^{\mathsf{T}})$  for all  $S \in \mathbb{R}^{m \times n}$ , and  $C_{pm}(S \otimes T) = (T \otimes S)C_{qn}$  for all  $T \in \mathbb{R}^{p \times q}$ . We will also make use of the property of the Kronecker product that  $(S \otimes T)^{-1} = S^{-1} \otimes T^{-1}$ .

It was shown in [26] that  $J^{T}J$  is rank-deficient, which necessitates the use of Levenberg-Marquardt by adding a proximal term. Denote

 $\gamma = \rho + \lambda_t$ , our task here is to efficiently invert the matrix  $J^T J + \gamma I$ . Notice that  $\gamma I = \gamma I \otimes I$  with appropriate sizes,  $(J^T J + \gamma I)^{-1}$  is equal to

$$\begin{bmatrix} (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} + \boldsymbol{\gamma}\boldsymbol{I}_{k}) \otimes \boldsymbol{I}_{m} & (\boldsymbol{I}_{k} \otimes \boldsymbol{W})\boldsymbol{C}_{kk}(\boldsymbol{I}_{k} \otimes \boldsymbol{H})^{\mathsf{T}} \\ (\boldsymbol{I}_{k} \otimes \boldsymbol{H})\boldsymbol{C}_{kk}(\boldsymbol{I}_{k} \otimes \boldsymbol{W})^{\mathsf{T}} & (\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} + \boldsymbol{\gamma}\boldsymbol{I}_{k}) \otimes \boldsymbol{I}_{n} \end{bmatrix}^{-1}.$$
 (7)

We now invoke the block-wise inversion formula

$$\begin{bmatrix} A_1 & A_2 \\ A_2^{\mathsf{T}} & A_3 \end{bmatrix}^{-1} =$$

$$\begin{bmatrix} (A_1 - A_2 A_3^{-1} A_2^{\mathsf{T}})^{-1} & -A_1^{-1} A_2 (A_3 - A_2^{\mathsf{T}} A_1^{-1} A_2)^{-1} \\ -A_3^{-1} A_2^{\mathsf{T}} (A_1 - A_2 A_3^{-1} A_2^{\mathsf{T}})^{-1} & (A_3 - A_2^{\mathsf{T}} A_1^{-1} A_2)^{-1} \end{bmatrix}.$$
(8)

The upper-left block of (7) equals to

$$\begin{pmatrix} (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} + \gamma \boldsymbol{I}_{k}) \otimes \boldsymbol{I}_{m} - \\ (\boldsymbol{I}_{k} \otimes \boldsymbol{W}) \begin{pmatrix} \boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} \otimes (\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} + \gamma \boldsymbol{I})^{-1} \end{pmatrix} (\boldsymbol{I}_{k} \otimes \boldsymbol{W})^{\mathsf{T}} \end{pmatrix}^{-1}.$$

Then we invoke the matrix inversion lemma to further simplify the upper-left block of (7) as

$$(\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} + \gamma \boldsymbol{I}_{k})^{-1} \otimes \boldsymbol{I}_{m} + \left( (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} + \gamma \boldsymbol{I}_{k})^{-1} \otimes \boldsymbol{W} \right) \times$$
(9a)

$$\left( (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H})^{-1} \otimes (\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} + \gamma \boldsymbol{I}) - (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} + \gamma \boldsymbol{I})^{-1} \otimes \boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} \right)^{-1}$$
(9b)

$$\times \left( \left( \boldsymbol{H}^{\mathsf{T}} \boldsymbol{H} + \gamma \boldsymbol{I}_{k} \right)^{-1} \otimes \boldsymbol{W} \right)^{\mathsf{T}}.$$
 (9c)

As we can see, this matrix is highly structured with the use of Kronecker product and a small  $k \times k$  matrix inversion  $(\mathbf{H}^{\mathsf{T}}\mathbf{H} + \gamma \mathbf{I}_k)^{-1}$ , except for the matrix inversion in (9b), which is a  $k^2 \times k^2$  matrix.

It turns out that (9b) can stil be efficiently calculated with complexity  $\mathcal{O}(k^3)$  by using the eigen-decompositions of  $W^{\top}W$  and  $H^{\top}H$ . Denote their eigen-decompositions as

$$\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} = \boldsymbol{Q}_{W}\boldsymbol{\Lambda}_{W}\boldsymbol{Q}_{W}^{\mathsf{T}}, \quad \boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} = \boldsymbol{Q}_{H}\boldsymbol{\Lambda}_{H}\boldsymbol{Q}_{H}^{\mathsf{T}}, \tag{10}$$

then (9b) equals to

$$(\boldsymbol{\mathcal{Q}}_{H}\otimes\boldsymbol{\mathcal{Q}}_{W})\left(\boldsymbol{\Lambda}_{H}^{-1}\otimes(\boldsymbol{\Lambda}_{W}+\gamma\boldsymbol{I})-(\boldsymbol{\Lambda}_{H}+\gamma\boldsymbol{I})^{-1}\otimes\boldsymbol{\Lambda}_{W}\right)^{-1}(\boldsymbol{\mathcal{Q}}_{H}\otimes\boldsymbol{\mathcal{Q}}_{W})^{\mathsf{T}}.$$

The matrix inverse in the middle is a  $k^2 \times k^2$  diagonal matrix, so the overall complexity is dominated by the eigen-decompostion of two  $k \times k$  matrices (10) with  $O(k^3)$  flops. In fact forming these matrices takes  $O(mk^2)$  and  $O(nk^2)$  flops respectively, so that the computations of eigen-decompositions are amortized.

Due to symmetry, the lower-right block of (7) can be similarly calculated by inter-changing W and H in (9). Finally, according the block-wise inversion formula, the off-diagonal blocks are simply some additional matrix multiplications, again with nice Kronecker structures and small-sized matrix inversions  $(W^TW + \gamma I)^{-1}$  and  $(H^TH + \gamma I)^{-1}$ .

The reader may notice that we never counted the computational complexity of calculating the Kronecker product of two matrices. This is because the Kronecker product is never explicitly calculated in our implementation—in the sequel we will see that by using another useful property of the Kronecker product  $(B^T \otimes A) \operatorname{vec}(S) = \operatorname{vec}(ASB)$ , we manage to reduce the per-iteration complexity of this Gauss-Newton-like method down to O(mnk).

## C. NLS-ADMM for NMF

Now we apply the previous result to the context of NLS-ADMM for NMF, with key computations shown in (5). Let us first calculate  $J_t^{\mathsf{T}} f(z_t)$ , which equals to

$$\begin{aligned} \boldsymbol{J}_{t}^{\mathsf{T}}\boldsymbol{f}(\boldsymbol{z}_{t}) &= \begin{bmatrix} \boldsymbol{H}_{t}^{\mathsf{T}} \otimes \boldsymbol{I}_{m} \\ (\boldsymbol{W}_{t}^{\mathsf{T}} \otimes \boldsymbol{I}_{n})\boldsymbol{C}_{mn} \end{bmatrix} \operatorname{vec}(\boldsymbol{W}_{t}\boldsymbol{H}_{t}^{\mathsf{T}} - \boldsymbol{X}) \\ &= \begin{bmatrix} \operatorname{vec}(\boldsymbol{W}_{t}\boldsymbol{H}_{t}^{\mathsf{T}}\boldsymbol{H}_{t} - \boldsymbol{X}\boldsymbol{H}_{t}) \\ \operatorname{vec}(\boldsymbol{H}_{t}\boldsymbol{W}_{t}^{\mathsf{T}}\boldsymbol{W}_{t} - \boldsymbol{X}^{\mathsf{T}}\boldsymbol{W}_{t}) \end{bmatrix}. \end{aligned}$$

Partition the dual variable  $\boldsymbol{u}$  in (6) as  $[\operatorname{vec}(\boldsymbol{U}_W)^\top \operatorname{vec}(\boldsymbol{U}_H)^\top]^\top$ with  $\boldsymbol{U}_W \in \mathbb{R}^{m \times k}$  and  $\boldsymbol{U}_H \in \mathbb{R}^{n \times k}$ . Define

$$W = W_t H_t^{\dagger} H_t - X H_t + \rho(W_t - W + U_W),$$
  
$$\widehat{H} = H_t W_t^{\dagger} W_t - X^{\dagger} W_t + \rho(H_t - H + U_H),$$

then

$$\boldsymbol{J}_t^{\mathsf{T}} \boldsymbol{f}(\boldsymbol{z}_t) + \rho(\boldsymbol{z}_t - \boldsymbol{z} + \boldsymbol{u}) = \begin{bmatrix} \operatorname{vec}(\boldsymbol{W}) \\ \operatorname{vec}(\boldsymbol{\widehat{H}}) \end{bmatrix}$$

Partition  $(\boldsymbol{J}_t^{\mathsf{T}} \boldsymbol{J}_t + \gamma \boldsymbol{I})^{-1}$  into four blocks with appropriate sizes, then

$$\left( \boldsymbol{J}_{t}^{\mathsf{T}} \boldsymbol{J}_{t} + \gamma \boldsymbol{I} \right)^{-1} \left( \boldsymbol{J}_{t}^{\mathsf{T}} \boldsymbol{f}(\boldsymbol{z}_{t}) + \rho(\boldsymbol{z}_{t} - \boldsymbol{z} + \boldsymbol{u}) \right) = \begin{bmatrix} \boldsymbol{G}_{1} & \boldsymbol{G}_{2} \\ \boldsymbol{G}_{2}^{\mathsf{T}} & \boldsymbol{G}_{3} \end{bmatrix} \begin{bmatrix} \operatorname{vec}(\widehat{\boldsymbol{W}}) \\ \operatorname{vec}(\widehat{\boldsymbol{H}}) \end{bmatrix}$$

We start by calculating  $G_1 \operatorname{vec}(\widehat{W})$ , where  $G_1$  is equal to (9). First

$$\left( (\boldsymbol{H}_t^{\mathsf{T}} \boldsymbol{H}_t + \gamma \boldsymbol{I}_k)^{-1} \otimes \boldsymbol{I}_m \right) \operatorname{vec}(\widehat{\boldsymbol{W}}) = \operatorname{vec}\left( \widehat{\boldsymbol{W}} (\boldsymbol{H}_t^{\mathsf{T}} \boldsymbol{H}_t + \gamma \boldsymbol{I}_k)^{-1} \right), \\ \left( (\boldsymbol{H}_t^{\mathsf{T}} \boldsymbol{H}_t + \gamma \boldsymbol{I}_k)^{-1} \otimes \boldsymbol{W}_t \right)^{\mathsf{T}} \operatorname{vec}(\widehat{\boldsymbol{W}}) = \operatorname{vec}\left( \boldsymbol{W}_t^{\mathsf{T}} \widehat{\boldsymbol{W}} (\boldsymbol{H}_t^{\mathsf{T}} \boldsymbol{H}_t + \gamma \boldsymbol{I}_k)^{-1} \right).$$

Now let  $\mathbf{K}_W = \mathbf{W}_t^{\top} \mathbf{\widehat{W}} (\mathbf{H}_t^{\top} \mathbf{H}_t + \gamma \mathbf{I}_k)^{-1}$ , and define a  $k \times k$  matrix S with its (i, j)-th value equal to

$$S(i,j) = \left(\frac{\lambda_W(i) + \gamma}{\lambda_H(j)} - \frac{\lambda_W(i)}{\lambda_H(j) + \gamma}\right)^{-1}$$

where  $\lambda_W(i)$  denotes the *i*-th eigenvalue of  $W^{\top}W$  and  $\lambda_H(j)$  denotes the *j*-th eigenvalue of  $H^{\top}H$ . Then

$$\left( (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H})^{-1} \otimes (\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} + \gamma \boldsymbol{I}) - (\boldsymbol{H}^{\mathsf{T}}\boldsymbol{H} + \gamma \boldsymbol{I})^{-1} \otimes \boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} \right)^{-1} \operatorname{vec}(\boldsymbol{K}_{W}) = \operatorname{vec} \left( \boldsymbol{\mathcal{Q}}_{W} \left( \boldsymbol{S} * \boldsymbol{\mathcal{Q}}_{W}^{\mathsf{T}}\boldsymbol{K}_{W} \boldsymbol{\mathcal{Q}}_{H} \right) \boldsymbol{\mathcal{Q}}_{H}^{\mathsf{T}} \right),$$

where \* denotes matrix Hadamard product, i.e., element-wise multiplication. This part is then multiplied by  $(\mathbf{H}^T\mathbf{H} + \gamma \mathbf{I}_k)^{-1} \otimes \mathbf{W}$ , and we obtain  $G_1 \operatorname{vec}(\widehat{\mathbf{W}})$ . Reverse the role of  $\mathbf{W}$  and  $\mathbf{H}$ ,  $G_3 \operatorname{vec}(\widehat{\mathbf{H}})$  can be similarly calculated.

Finally, let  $\overline{H} = G_3 \operatorname{vec}(\widehat{H})$ , we have

$$\begin{split} & \boldsymbol{G}_{2}^{\mathsf{T}}\operatorname{vec}(\widehat{\boldsymbol{W}}) \\ &= -\left( (\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} + \gamma \boldsymbol{I}_{k})^{-1} \otimes \boldsymbol{I}_{n} \right) (\boldsymbol{I}_{k} \otimes \boldsymbol{W}) \boldsymbol{C}_{kk} (\boldsymbol{I}_{k} \otimes \boldsymbol{H})^{\mathsf{T}}\operatorname{vec}(\overline{\boldsymbol{H}}) \\ &= -\operatorname{vec}\left( \boldsymbol{W}(\overline{\boldsymbol{H}}^{\mathsf{T}}\boldsymbol{H}) (\boldsymbol{W}^{\mathsf{T}}\boldsymbol{W} + \gamma \boldsymbol{I}_{k})^{-1} \right). \end{split}$$

Similarly, one can calculate  $G_2 \operatorname{vec}(\widehat{H})$  by reversing the role of W and H, and

$$\left(\boldsymbol{J}_{\boldsymbol{I}}^{\mathsf{T}}\boldsymbol{J}_{\boldsymbol{I}}+\boldsymbol{\gamma}\boldsymbol{I}\right)^{-1}\begin{bmatrix}\operatorname{vec}(\widehat{\boldsymbol{W}})\\\operatorname{vec}(\widehat{\boldsymbol{H}})\end{bmatrix}=\begin{bmatrix}\boldsymbol{G}_{1}\operatorname{vec}(\widehat{\boldsymbol{W}})+\boldsymbol{G}_{2}\operatorname{vec}(\widehat{\boldsymbol{H}})\\\boldsymbol{G}_{3}\operatorname{vec}(\widehat{\boldsymbol{H}})+\boldsymbol{G}_{2}^{\mathsf{T}}\operatorname{vec}(\widehat{\boldsymbol{W}})\end{bmatrix}$$

The overall algorithm is summarized in Algorithm 1.

## Algorithm 1 NLS-ADMM for NMF

1: initialize W and H2:  $\rho = 1$ 3: repeat  $W^{\mathsf{T}}W = Q_W \Lambda_W Q_W^{\mathsf{T}}$ 4: 5: ▷ eigen-decomposition 6:  $\widetilde{\boldsymbol{Z}}_{W} = \boldsymbol{W}, \widetilde{\boldsymbol{Z}}_{H} = \boldsymbol{H}, \boldsymbol{U}_{W} = \boldsymbol{U}_{H} = 0$ 7: ▷ initialize ADMM repeat 8:  $\widehat{\boldsymbol{W}} \leftarrow \boldsymbol{W} \boldsymbol{H}^{\mathsf{T}} \boldsymbol{H} - \boldsymbol{F}_{\boldsymbol{W}} + \rho(\boldsymbol{W} - \boldsymbol{Z}_{\boldsymbol{W}} + \boldsymbol{U}_{\boldsymbol{W}})$ 9:  $\widehat{\boldsymbol{H}} \leftarrow \boldsymbol{H} \boldsymbol{W}^{\!\!\top} \boldsymbol{W} - \boldsymbol{F}_{H} + \rho(\boldsymbol{H} - \boldsymbol{Z}_{H} + \boldsymbol{U}_{H})$  $K_W \leftarrow W^{\mathsf{T}} \widehat{W} (H^{\mathsf{T}} H + \gamma I_k)^{-1}$ 10:  $\boldsymbol{K}_{\boldsymbol{H}} \leftarrow \boldsymbol{H}^{\mathsf{T}} \widehat{\boldsymbol{H}} (\boldsymbol{W}^{\mathsf{T}} \boldsymbol{W} + \gamma \boldsymbol{I}_{k})^{-1}$  $K_W \leftarrow Q_W \left( S * Q_W^\top K_W Q_H \right) Q_H^\top$ 11:  $K_H \leftarrow Q_H ((1/S) * Q_H^{\mathsf{T}} K_H Q_W) Q_W^{\mathsf{T}}$  $\overline{W} \leftarrow \widehat{W} (H^{\top}H + \gamma I_k)^{-1} + WK_W (H^{\top}H + \gamma I_k)^{-1} \\ \overline{H} \leftarrow \widehat{H} (W^{\top}W + \gamma I_k)^{-1} + HK_H (W^{\top}W + \gamma I_k)^{-1}$ 12:  $\widetilde{Z}_{W} \leftarrow W - \overline{W} + H(\overline{W}^{\mathsf{T}}W)(H^{\mathsf{T}}H + \gamma I_{k})^{-1}$  $\widetilde{Z}_{H} \leftarrow H - \overline{H} + W(\overline{H}^{\mathsf{T}}H)(W^{\mathsf{T}}W + \gamma I_{k})^{-1}$ 13:  $Z_W \leftarrow [\widetilde{Z}_W + U_W]_+$ 14:  $Z_H \leftarrow [\widetilde{Z}_H + U_H]_+$  $U_W \leftarrow U_W + \widetilde{Z}_W - Z_W$ 15:  $U_H \leftarrow U_H + \widetilde{Z}_H - Z_H$ until convergence 16: if  $||X - WH^{\top}||_{F}^{2} < ||X - Z_{W}Z_{H}^{\top}||_{F}^{2}$  then 17:  $\triangleright$  re-calculate the update 18:  $\rho \leftarrow 2\rho$ 19: else  $W \leftarrow Z_W$ 20:  $H \leftarrow Z_H$  $\triangleright$  accept the update 21:  $\rho \leftarrow \rho/2$ 22: end if 23: until convergence

## **III. NUMERICAL EXPERIMENTS**

We conclude the paper by conducting some experiment on synthetic data to showcase the fast convergence rate of the proposed NLS-ADMM for NMF. All experiments are conducted in MATLAB 2018b on an iMac Pro with 10 Xeon processers and 32GB memory. We compare with two baseline algorithms, AO-ADMM [21] and APG with extrapolation [18], both being reported to work the best in a lot of cases.

We let m = 100, n = 150, k = 10, and randomly generate the ground-truth factors W and H from uniform distribution between [0, 1]. Then the data matrix  $X = WH^{\top}$  is constructed. Although there is no noise, it is known that finding an exact factorization is relatively hard because the factors are not identifiable (since they are all dense) [12].

A typical convergence plot is shown in Figure 1. Indeed, even though we have greatly reduced the per-iteration complexity of NLS-ADMM, every iteration still takes longer time than that of AO-ADMM and APG; but thanks to the exploitation of second-order information in the loss function, NLS-ADMM is able to escape the swamp effect, which often happens in alternating-type algorithms. The result is not surprising as it resembles the cases people have observed in CPD when alternating optimization is compared with quasi-Newton methods—alternating least-squares (ALS) sometimes gets stucked at a non-critical point for a very long time, whereas Gauss-Newton or Levenberg-Marquardt is able to quickly converge



Fig. 1: One instance of convergence comparison on synthetic data.

TABLE I: Averaged performance of NMF algorithms

Algorithms	$\ \boldsymbol{X} - \boldsymbol{W}\boldsymbol{H}^{T}\ _{F}^{2}$	run time	iterations
NLS-ADMM	$2.18 \times 10^{-8}$	1.1375	23.23
AO-ADMM	0.0108	7.4724	42019
APG	0.0099	10.9181	65035

to a global minimum with super-linear convegence rate [24]. It is interesting to see that by using ADMM to solve the constrained nonlinear least-squares problem, such good convegence property still exists.

The above experiment setting is randomly repeated 100 times and the averaged performances are shown in Table I. As we can see again, NLS-ADMM is able to obtain remarkably accurate solution with shortest amount of run time and very small number of iterations.

## IV. CONCLUSION

This paper serves as the first attempt to apply the algorithmic framework of Gauss-Newton / Levenberg-Marquardt for the widely used nonnegative matrix factorization (NMF) problem. There are two notable contributions when applying this seemingly simple idea. The first resolves the challenge of how to properly handle the nonnegativity constraints on the latent factors, which is an important feature of NMF, but cannot be handled by the original Gauss-Newton framework. Our solution is to solve a constrained leastsquares instead in each iteration, using alternating direction method of multipliers (ADMM). Another challenge is the high per-iteration complexity to naively solve a (constrained) least-squares problem due to the "all-at-once" update nature of Gauss-Newton. By carefully exploiting the structure of the Jacobian Grammian, we managed to bring down the per-iteration complexity to the same order of that of an alternating optimization algorithm such as AO-ADMM and APG. Numerical experiments on synthetic data shows the superior convergence rate compared to the state-of-the-arts.

For future work, we notice that it is still taking many ADMM iterations for each GN updates, especially at the beginning stage. An even faster algorithm to solve the constrained least-squares subproblem would greatly accelerate the proposed proximal Gauss-Newton algorithmic framework. It is also possible to extend the framework for constrained tensor factorization. However, to the best of our knowledge, there has not been a direct method of inverting the Jacobian Grammian matrix with complexity similar to an ALS step. A more in-depth study is required to make it work for tensors.

#### REFERENCES

- X. Fu, K. Huang, N. D. Sidiropoulos, and W.-K. Ma, "Nonnegative matrix factorization for signal and data analytics: Identifiability, algorithms, and applications," *IEEE Signal Processing Magazine*, vol. 36, pp. 59–80, 2019.
- [2] D. D. Lee and H. S. Seung, "Learning the parts of objects by nonnegative matrix factorization," *Nature*, vol. 401, no. 6755, p. 788, 1999.
- [3] P. Smaragdis, C. Fevotte, G. J. Mysore, N. Mohammadiha, and M. Hoffman, "Static and dynamic source separation using nonnegative factorizations: A unified view," *IEEE Signal Processing Magazine*, vol. 31, no. 3, pp. 66–75, 2014.
- [4] X. Fu, W.-K. Ma, K. Huang, and N. D. Sidiropoulos, "Blind separation of quasi-stationary sources: Exploiting convex geometry in covariance domain," *IEEE Transactions on Signal Processing*, vol. 63, no. 9, pp. 2306–2320, 2015.
- [5] W.-K. Ma, J. M. Bioucas-Dias, T.-H. Chan, N. Gillis, P. Gader, A. J. Plaza, A. Ambikapathi, and C.-Y. Chi, "A signal processing perspective on hyperspectral unmixing: Insights from remote sensing," *IEEE Signal Processing Magazine*, vol. 31, no. 1, pp. 67–81, 2014.
- [6] X. Fu, K. Huang, B. Yang, W.-K. Ma, and N. D. Sidiropoulos, "Robust volume minimization-based matrix factorization for remote sensing and document clustering," *IEEE Transactions on Signal Processing*, vol. 64, no. 23, pp. 6254–6268, 2016.
- [7] S. Arora, R. Ge, Y. Halpern, D. Mimno, A. Moitra, D. Sontag, Y. Wu, and M. Zhu, "A practical algorithm for topic modeling with provable guarantees," in *International Conference on Machine Learning*, 2013, pp. 280–288.
- [8] K. Huang, X. Fu, and N. D. Sidiropoulos, "Anchor-free correlated topic modeling: Identifiability and algorithm," in Advances in Neural Information Processing Systems, 2016, pp. 1786–1794.
- [9] K. Huang, X. Fu, and N. Sidiropoulos, "Learning hidden Markov models from pairwise co-occurrences with application to topic modeling," in *International Conference on Machine Learning*, 2018, pp. 2073–2082.
- [10] J. Yang and J. Leskovec, "Overlapping community detection at scale: A nonnegative matrix factorization approach," in ACM International Conference on Web Search and Data Mining. ACM, 2013, pp. 587–596.
- [11] K. Huang and X. Fu, "Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm," in *International Conference on Machine Learning*, 2019, pp. 2859–2868.
- [12] K. Huang, N. D. Sidiropoulos, and A. Swami, "Non-negative matrix factorization revisited: Uniqueness and algorithm for symmetric decomposition," *IEEE Transactions on Signal Processing*, vol. 62, no. 1, pp. 211–224, 2013.
- [13] X. Fu, K. Huang, and N. D. Sidiropoulos, "On identifiability of nonnegative matrix factorization," *IEEE Signal Processing Letters*, vol. 25, no. 3, pp. 328–332, 2018.
- [14] S. A. Vavasis, "On the complexity of nonnegative matrix factorization," *SIAM Journal on Optimization*, vol. 20, no. 3, pp. 1364–1377, 2009.
- [15] D. D. Lee and H. S. Seung, "Algorithms for non-negative matrix factorization," in Advances in Neural Information Processing Systems, 2001, pp. 556–562.
- [16] A. Cichocki and A.-H. Phan, "Fast local algorithms for large scale nonnegative matrix and tensor factorizations," *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, vol. 92, no. 3, pp. 708–721, 2009.
- [17] C.-J. Lin, "Projected gradient methods for nonnegative matrix factorization," *Neural Computation*, vol. 19, no. 10, pp. 2756–2779, 2007.
- [18] Y. Xu and W. Yin, "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion," *SIAM Journal on Imaging Sciences*, vol. 6, no. 3, pp. 1758–1789, 2013.
- [19] H. Kim and H. Park, "Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method," *SIAM Journal on Matrix Analysis and Applications*, vol. 30, no. 2, pp. 713–730, 2008.
- [20] J. Kim and H. Park, "Fast nonnegative matrix factorization: An activeset-like method and comparisons," *SIAM Journal on Scientific Computing*, vol. 33, no. 6, pp. 3261–3281, 2011.
- [21] K. Huang, N. Sidiropoulos, and A. P. Liavas, "A flexible and efficient algorithmic framework for constrained matrix and tensor factorization," *IEEE Transactions on Signal Processing*, vol. 64, no. 19, pp. 5052–5065, 2016.

- [22] Y. Xu, W. Yin, Z. Wen, and Y. Zhang, "An alternating direction algorithm for matrix completion with nonnegative factors," *Frontiers of Mathematics in China*, vol. 7, no. 2, pp. 365–384, 2012.
- [23] D. Hajinezhad, T.-H. Chang, X. Wang, Q. Shi, and M. Hong, "Nonnegative matrix factorization using ADMM: Algorithm and convergence analysis," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 4742–4746.
- [24] L. Sorber, M. Van Barel, and L. De Lathauwer, "Optimization-based algorithms for tensor decompositions: Canonical polyadic decomposition, decomposition in rank- $(L_r, L_r, 1)$  terms, and a new generalization," *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 695–720, 2013.
- [25] A.-H. Phan, P. Tichavsky, and A. Cichocki, "Low complexity damped Gauss-Newton algorithms for CANDECOMP/PARAFAC," *SIAM Journal on Matrix Analysis and Applications*, vol. 34, no. 1, pp. 126–147, 2013.
- [26] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos, "Tensor decomposition for signal processing and machine learning," *IEEE Transactions on Signal Processing*, vol. 65, no. 13, pp. 3551–3582.
- [27] S. Boyd and L. Vandenberghe, Introduction to Applied Linear Algebra: Vectors, Matrices, and Least Squares. Cambridge University Press, 2018.
- [28] S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein *et al.*, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends*® *in Machine learning*, vol. 3, no. 1, pp. 1–122, 2011.