

NIH Public Access

Author Manuscript

J Signal Process Syst. Author manuscript; available in PMC 2013 June 05.

Published in final edited form as:

J Signal Process Syst. 2011 January 12; 65(3): 351-359. doi:10.1007/s11265-010-0511-8.

Correlated Noise: How it Breaks NMF, and What to Do About It

Sergey M. Plis,

Computer Science Department, University of New Mexico, New Mexico, USA, 87131

Vamsi K. Potluru,

Computer Science Department, University of New Mexico, New Mexico, USA, 87131

Terran Lane, and

Computer Science Department, University of New Mexico, New Mexico, USA, 87131

Vince D. Calhoun

Electrical and Computer Engineering Department, University of New Mexico, New Mexico, USA, 87131, Mind Research Network, New Mexico, USA, 87131

Abstract

Non-negative matrix factorization (NMF) is a problem of decomposing multivariate data into a set of features and their corresponding activations. When applied to experimental data, NMF has to cope with noise, which is often highly correlated. We show that correlated noise can break the Donoho and Stodden separability conditions of a dataset and a regular NMF algorithm will fail to decompose it, even when given freedom to be able to represent the noise as a separate feature. To cope with this issue, we present an algorithm for NMF with a generalized least squares objective function (glsNMF) and derive multiplicative updates for the method together with proving their convergence. The new algorithm successfully recovers the true representation from the noisy data. Robust performance can make glsNMF a valuable tool for analyzing empirical data.

1 Introduction

Since the introduction of multiplicative updates for non-negative matrix factorization (NMF) [1], the algorithm has gained general recognition. Simplicity of implementation, an adaptive learning rate and automatic satisfaction of positivity constraints are in part responsible for the wide acceptance of the algorithm. It has been successfully used to analyze functional brain imaging data [2–4], gene expression [5], and other empirical datasets.

Lee and Seung [1] provide two updates for NMF: one is based on the least squares (LS) criteria and the other on Kullback-Leibler (KL) divergence. In this study we focus on LS updates, for which the data model is:

 $X = WH + \varepsilon$, (1)

where each entry in X, W and H is greater than or equal to zero and e is Gaussian noise. Subsequent sections provide further details.

The LS formulation implicitly assumes that the noise is white. This is a widely used assumption and it is valid in many realistic cases with a large number of independent noise sources. However, in many experimental settings noise is more complicated and is not limited to white sensor noise. In these environments, noise represents background activity, which can have complex covariance structure. Ignoring the structure in the noise can change the results of NMF substantially.

Donoho and Stodden introduced the notion of a separable factorial articulation family [6] as a collection of points obeying three conditions: generative model, separability and complete factorial sampling. Datasets satisfying these conditions are guaranteed to be properly factored by any correct NMF algorithm. The presence of correlated noise may, however, violate the conditions and render otherwise separable dataset not factorable by NMF. We show an example in which a dataset that otherwise satisfies the Donoho and Stodden conditions is not factored properly by the regular NMF when contaminated with correlated noise. We also show that despite a reasonable expectation that an NMF model of sufficient rank can recover correlated components in the noise simply as nuisance features, NMF fails to do so and most features are recovered incorrectly. Undermodeling the noise may lead to a misinterpretation of results when applied to real dataset without known ground truth.

As a solution, we introduce a generalized least squares NMF (glsNMF) algorithm that takes the correlation structure of the noise into account. We derive multiplicative updates for the algorithm providing a proof of their convergence and demonstrate the algorithm's performance on a synthetic dataset. We show that the new algorithm handles the noisy data better than the LS based algorithm and produces the expected unique factors.

2 NMF

NMF is a tool producing a low rank approximation to a non-negative data matrix by splitting it into a product of two non-negative matrix factors. The constraint of nonnegativity (all elements are 0) usually results in a parts-based representation making NMF different from other factorization techniques yielding more holistic representations, such as principal component analysis (PCA) and vector quantization (VQ).

Using standard notation [1,6], we formally define NMF task as follows. Given a nonnegative $m \times n$ matrix X, represent it with a product of two non-negative matrices W, H of sizes $m \times r$ and $r \times n$ respectively:

X pprox WH. (2)

The non-negativity constraint corresponds to the intuitive notion of features adding up, without canceling each other, to give the resulting data.

Lee and Seung [1] describe two multiplicative update rules for W and H which work well in practice. The updates correspond to two different cost functions representing the quality of approximation. In this work, we use the Frobenius norm for the cost function:

$$E = \frac{1}{2} \left\| \boldsymbol{X} - \boldsymbol{W} \boldsymbol{H} \right\|_{F} \quad (3)$$

and the corresponding updates are:

$$W = W \odot \frac{XH^T}{WHH^T}$$
 (4)

$$H = H \odot \frac{W^T X}{W^T W H}, \quad (5)$$

where $\|.\|_F$ denotes the Frobenius norm and the operator \odot represents element-wise multiplication. Division is also element-wise. We have omitted iteration indices for clarity. It should be noted that the cost function to be minimized is convex in either *W* or *H* but not

in both [1]. As the algorithm iterates using the updates (4) and (5), W and H converge to a local minimum of the cost function.

The slightly mysterious form for the above updates can be understood as follows. Simple additive updates for H and W are given by:

$$H = H + \eta_H \odot (W^T X - W^T W H) \quad (6)$$
$$W = W + \eta_W \odot (X H^T - W H H^T) \quad (7)$$

If the learning rate given by the matrix elements of $\eta_{\{W,H\}}$ are all set to some small positive number then this is the conventional gradient descent. However, setting the learning rate matrices as follows:

$$\eta_{H} = \frac{H}{W^{T}WH}, \quad (8)$$
$$\eta_{W} = \frac{W}{WHH^{T}}, \quad (9)$$

w 1111

where division is element-wise, produces the NMF updates.

3 Failure Mode

An example of a separable factorial articulation family is the swimmer dataset presented in [6]. The dataset contains 256 32×32 images of a swimmer with all possible combinations of limbs positions (encoded in the feature or source matrix **W** with a column per 32×32 limb, and corresponding feature activations encoded in **H**), as shown in Figure 1.

In order to study the effect of the correlated noise on the algorithm we have constructed such noise, where a small part of the image is spatially correlated. Figure 2 shows several random samples from the swimmer dataset of Figure 1 contaminated by the correlated noise.

The LS objective function of (3) can be derived from the Gaussian likelihood with noise covariance of the form $\sigma^2 I$. Note that correlated noise results in a differently structured covariance matrix. The covariance of the correlated noise is shown in Figure 3. It is clearly very close to a diagonal matrix. For comparison, the figure also shows a close up image of a section of the covariance, where there are correlations. Correlations among the 2% of the image pixels are high as demonstrated by the high contrast of the covariance image in the figure. Several samples of the noise are shown in Figure 4. The correlated part has the shape of the swimmer's torso shifted to the left of the original torso position. In summary the noise is mostly white, with a small locally concentrated correlated component.

A reasonable expectation of NMF's behavior on this dataset would be a result that has the correlation torso as a separate feature with the other features correctly recovered. This common-sense behavior would go along with other matrix factorization techniques such as independent component analysis (ICA), which exploit this feature for denoising. Surprisingly, we have observed quite different behavior. A typical result is shown in Figure 5, where it becomes clear that correlated noise affects the estimation of all of the features, in addition to being estimated as a separate feature. For comparison the correct factorization that is obtained by NMF from the noiseless dataset is shown in Figure 6. Note that we have

observed similar behavior even when using the KL divergence based NMF objective, although we do not pursue this finding further in this work.

Introduction of the torso-shaped correlation in the noise violates the *separability* condition from [6]. The condition requires that each part/articulation pair's presence or absence in the image is indicated by a certain pixel associated to that pair. However, the torso-shaped correlation, when present, can overlap with limbs in some positions. Note that conditions of generative model and complete factorial sampling are still satisfied, because the correlation can always be treated as a separate feature.

4 Proposed Solution

We argue that in practical applications of NMF one needs to model the noise adequately. Here we propose an NMF algorithm that alleviates the problem caused by the correlated noise.

One of the objective functions for non-negative matrix factorization proposed in [1] is the least squares error (LSE) of (3). After rewriting (3) in the matrix form:

$$E = \frac{1}{2} \operatorname{Tr}((\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H})^{T}(\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H})), \quad (10)$$

the assumption of zero mean noise with unit standard deviation becomes explicit.

For optimization purposes, the formulation is also valid for noise with covariance structure $\sigma^2 I$. Richer noise structures, including those with diagonal covariance or correlated noise, are not captured by such a formulation. The former problem has been addressed in [7]. In this case the scaling of each dimension by a corresponding positive variance is performed. Scaling by the positive constants does not alter the original formulation of multiplicative updates.

We address the problem of generalized least squares (GLS) of (11), where C is the noise covariance, and derive multiplicative updates for this general form of the objective:

$$E = \frac{1}{2} \operatorname{Tr}((\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H})^T \boldsymbol{C}^{-1}(\boldsymbol{X} - \boldsymbol{W}\boldsymbol{H})) \quad (11)$$

4.1 Derivation of the updates

To simplify the notation, we define the precision matrix, $S = C^{-1}$. First, we rewrite the objective:

$$E = \frac{1}{2} (\operatorname{Tr}(\boldsymbol{X}^T \boldsymbol{S} \boldsymbol{X}) + \operatorname{Tr}(\boldsymbol{H}^T \boldsymbol{W}^T \boldsymbol{S} \boldsymbol{W} \boldsymbol{H}) - \operatorname{Tr}(\boldsymbol{X}^T \boldsymbol{S} \boldsymbol{W} \boldsymbol{H}) - \operatorname{Tr}(\boldsymbol{H}^T \boldsymbol{W}^T \boldsymbol{S} \boldsymbol{X})).$$
(12)

Then find the derivatives:

$$\frac{\partial E}{\partial W} = S(WHH^T - XH^T) \quad (13)$$

$$\frac{\partial E}{\partial \boldsymbol{H}} = \boldsymbol{W}^T \boldsymbol{S} (\boldsymbol{W} \boldsymbol{H} - \boldsymbol{X}), \quad (14)$$

where we have used the fact that S is symmetric. Next, we follow the procedure similar to the one used to obtain the updates for the original NMF algorithm [1]. We set the values of the learning rate for both matrices W and H such that the updates become multiplicative. A detailed convergence proof that lead to these values is presented in Appendix. Finally, we arrive at the following multiplicative updates for the GLS error function:

$$W = W \odot \frac{S^{+}XH^{T} + S^{-}WHH^{T}}{S^{-}XH^{T} + S^{+}WHH^{T}}$$
(15)

$$H = H \odot \frac{W^T S^+ X + W^T S^- W H}{W^T S^- X + W^T S^+ W H}$$
(16)

In these updates, the precision matrix is split into two parts as follows:

$$\begin{split} S_{ij}^{+} = \begin{cases} S_{ij} & S_{ij} > 0, \\ S_{ij} + \lambda & i = j, \\ 0 & \text{otherwise}, \end{cases} \\ S_{ij}^{-} = \begin{cases} |S_{ij}| & S_{ij} < 0, \\ |S_{ij}| + \lambda & i = j, \\ 0 & \text{otherwise}. \end{cases} \end{split}$$

In the matrix representation the split of *S* is expressed as:

$$S^{+}=S^{+}+\lambda I$$
$$S^{-}=\widehat{S}^{-}+\lambda I$$
$$S=\widehat{S}^{+}-\widehat{S}^{-}$$

where λ is the minimal negative eigenvalue of \hat{S}^- or 0 if \hat{S}^- is positive semidefinite or empty. This ensures that S^- is positive semidefinite – a condition required for convergence properties of the updates. We defer the proof to the appendix. The appendix also provides details of implementing a projected gradient algorithm for glsNMF.

4.2 Complexity

Introduction of S^+ and S^- added complexity to the updates, namely, four additional matrix multiplications and two matrix summations. Repeating parts of expressions in the numerator and the denominator of equations (16) and (15) can be precomputed before each respective updates. After that only multiplications by parts of the precision matrix and summation of the result are required.

4.3 Covariance estimation

In order for glsNMF to function properly, a good estimate of the noise covariance is required. This is sometimes possible to obtain as a sample covariance of the background measurements of the system without the task of interest. When such background noise is available, it is possible to estimate the covariance using the sample covariance estimation [8]. This is the method we use in this paper. This is especially true in functional electromagnetic brain imaging (an area of increasing use of NMF [4]), where sampling rates allow collection of sufficient samples at least for spatial only analysis. Many models that use the covariance matrix, like glsNMF, assume that it is computed elsewhere.

5 Results

The glsNMF algorithm was applied to the previously introduces noisy swimmer dataset. As before, the number of features was set to 20.

The features obtained by glsNMF are shown in Figure 7. Compare with the features typically obtained by NMF in Figure 5. Note that we have run both algorithms many times changing the starting point. The starting points in the figures are the same.

NMF applied to the noisy dataset produces features spanning several swimmer configurations at once. Thus, it is hard to identify the original noise free image. In addition to that there is a spurious feature – the correlated part of the noise. It is not only present as a feature by itself but also contaminates some other features extracted by the algorithm.

Features extracted by glsNMF are sparse and completely recover the limbs. Furthermore, the correlated part of the noise is not present as a separate feature or as a part of any other features. Although some residual noise still remains even after convergence it does not destroy the parts based representation.

6 Discussion

We have shown that NMF can fail in the case of data contaminated by correlated noise. This is almost always unavoidable situation when dealing with experimental datasets. Issue of noise has been previously addressed in NMF algorithms by rescaling each of the estimates by the amount of noise variance (uncertainty) [7], or by using Gaussian process priors to smooth out the estimates [9]. Results similar to [7] can probably be achieved using the advances in research on weighted NMF [10]. An approach that has formulation similar to our suggested solution was presented in [11]. However, there the goal was not to target correlated noise and also the novelty of our formulation is the multiplicative updates and their convergence proof. In fact, a solution by a projected gradient [12] method is easily possible and we have also derived it for our method including the derivation in Appendix.

There are a multitude of extensions to NMF, such as additional sparsity, convolutive algorithm etc. [13–15]. We believe some of them can benefit from using GLS objective function.

7 Conclusions

A growing interest in application of NMF to experimental datasets [3,4,2,5] requires special handling of issues introduces by unavoidable presence of noise. We have demonstrated that the NMF algorithm can fail in the presence of correlated noise which can violate the separability assumption of unique factorization and degrade the results in applications such as feature extraction. We also proposed the glsNMF algorithm as a solution to the correlated noise problem which is able to recover features from data with correlated noise. For this we have derived a multiplicative update and presented a convergence proof. Our future work includes application of the method to functional brain imaging and gene expression datasets as well as extending the method to deal with large dimensionality of the input space which makes the covariance matrix hard to handle. It is also possible to perform glsNMF with simultaneous estimation of the covariance matrix which we also leave for future work.

References

1. Lee, Daniel D.; Seung, Sebastian H. Algorithms for non-negative matrix factorization. NIPS. 2000:556–562.

- 3. Potluru, VK.; Calhoun, VD. Group learning using contrast NMF: Application to functional and structural MRI of schizophrenia. Circuits and Systems, 2008. ISCAS 2008. IEEE International Symposium; May 2008; p. 1336-1339.
- Mørup M, Hansen LK, Arnfred SM. ERPWAVELAB a toolbox for multi-channel analysis of timefrequency transformed event related potentials. Journal of Neuroscience Methods. 2007; 161:361– 368. [PubMed: 17204335]
- Devarajan, Karthik. Nonnegative matrix factorization: An analytical and interpretive tool in computational biology. PLoS Comput Biol. Jul.2008 4(7):e1000029. [PubMed: 18654623]
- 6. Donoho, David; Stodden, Victoria. When does non-negative matrix factorization give a correct decomposition into parts?. In: Thrun, Sebastian; Saul, Lawrence; Schölkopf, Bernhard, editors. Advances in Neural Information Processing Systems. Vol. 16. MIT Press; Cambridge, MA: 2004.
- 7. Wang G, Kossenkov AV, Ochs MF. LS-NMF: a modified non-negative matrix factorization algorithm utilizing uncertainty estimates. BMC Bioinformatics. 2006; 7
- 8. Stark, H.; Woods, JW. Probability and random processes with applications to signal processing. Prentice Hall; Upper Saddle River, NJ: 2002.
- Schmidt, Mikkel N.; Laurberg, Hans. Nonnegative matrix factorization with Gaussian process priors. Computational intelligence and neuroscience. 2008:361705. [PubMed: 18464923]
- Guillamet D, Vitri J, Schiele B. Introducing a weighted non-negative matrix factorization for image classification. Pattern Recognition Letters. 2003; 24(14):2447–2454.
- Cichocki, Andrzej; Zdunek, Rafal. Regularized alternating least squares algorithms for nonnegative matrix/tensor factorization. ISNN 07: Proceedings of the 4th international symposium on Neural Networks; Berlin, Heidelberg. Springer-Verlag; 2007. p. 793-802.
- Lin, Chih-Jen. Projected gradient methods for nonnegative matrix factorization. Neural Comp. Oct; 2007 19(10):2756–2779.
- Hoyer, PO. Non-negative sparse coding. Neural Networks for Signal Processing, 2002. Proceedings of the 2002 12th IEEE Workshop on; 2002. p. 557-565.
- Potluru, VK.; Plis, SM.; Calhoun, VD. Sparse shift-invariant NMF. Image Analysis and Interpretation, 2008. SSIAI 2008. IEEE Southwest Symposium; mar 2008; p. 69-72.
- 15. O'Grady, Paul D.; Pearlmutter, Barak A. Convolutive non-negative matrix factorisation with a sparseness constraint. Proceedings of the IEEE International Workshop on Machine Learning for Signal Processing (MLSP 2006); Maynooth, Ireland. sep 2006; p. 427-432.
- 16. Petersen KB, Pedersen MS. The matrix cookbook. Oct.2008 Version 20081110.

Appendix

Convergence proof

Consider the problem for a single column of H denoted by h. The corresponding column of X is given by x. The objective is now given by:

$$F(\boldsymbol{h}) = \frac{1}{2} (\boldsymbol{x} - \boldsymbol{W}\boldsymbol{h})^T \boldsymbol{S} (\boldsymbol{x} - \boldsymbol{W}\boldsymbol{h}) \quad (17)$$

We define an auxiliary function G(h, v) with the properties that G(h, h) = F(h) and G(h, h')

F(h), where h^k is the current estimate with iteration index k and h is the free parameter. The multiplicative update rule is found at each iteration by minimizing the auxiliary function:

$$h^{k+1} = \arg\min_{h} G(h, h^k)$$
 (18)

Note that this does not increase the objective function F, as we have

$$F(\boldsymbol{h}^{k+1}) \leq G(\boldsymbol{h}^{k+1}, \boldsymbol{h}^k) \leq G(\boldsymbol{h}^k, \boldsymbol{h}^k) = F(\boldsymbol{h}^k) \quad (19)$$

Define G as follows:

$$G(\boldsymbol{h},\boldsymbol{h}^{k}) = F(\boldsymbol{h}^{k}) + (\boldsymbol{h} - \boldsymbol{h}^{k})^{T} \nabla F(\boldsymbol{h}^{k}) + \frac{1}{2} (\boldsymbol{h} - \boldsymbol{h}^{k})^{T} \boldsymbol{K}(\boldsymbol{h}^{k}) (\boldsymbol{h} - \boldsymbol{h}^{k})$$
(20)

where the diagonal matrix $K(h^k)$ is defined as

$$\boldsymbol{K}_{ab}(\boldsymbol{h}^{k}) = \delta_{ab} \frac{(\boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{W} \boldsymbol{h}^{k} + \boldsymbol{W}^{T} \boldsymbol{S}^{-} \boldsymbol{x})_{a}}{\boldsymbol{h}_{a}^{k}}, \quad (21)$$

where a and b are indices of vector and matrix elements, and

$$\delta_{ab} = \begin{cases} 1, & \text{if } a = b \\ 0, & \text{if } a \neq b \end{cases} (22)$$

It trivially holds that G(h, h) = F(h). For $G(h, h^k) = F(h)$ to hold we need

$$(\boldsymbol{h} - \boldsymbol{h}^k)^T [\boldsymbol{K}(\boldsymbol{h}^k) - \boldsymbol{W}^T \boldsymbol{S} \boldsymbol{W}](\boldsymbol{h} - \boldsymbol{h}^k) \ge 0$$
 (23)

We denote $M = K(h^k) - W^T SW$, observing that M must be positive semidefinite for (23) to hold, and split it into parts as follows:

$$M = P_1 + P_2 + P_3 \quad (24)$$

$$(\boldsymbol{P}_{1})_{ab} = \delta_{ab} \frac{(\boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{W} \boldsymbol{h}^{k})_{a}}{\boldsymbol{h}^{k}_{a}} - (\boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{W})_{ab} \quad (25)$$
$$(\boldsymbol{P}_{2})_{ab} = \delta_{ab} \frac{(\boldsymbol{W}^{T} \boldsymbol{S}^{-} \boldsymbol{x})_{a}}{\boldsymbol{h}^{k}_{a}} \quad (26)$$
$$(\boldsymbol{P}_{3})_{ab} = (\boldsymbol{W}^{T} \boldsymbol{S}^{-} \boldsymbol{W})_{ab} \quad (27)$$

If each P_i is positive semidefinite then their sum M is also so. The most difficult is P_1 and we show it is positive semidefinite the last. P_2 is trivially positive semidefinite since it is a diagonal matrix with non-negative entries. P_3 is also positive semidefinite since by construction in Section 4.1 we obtain a positive semidefinite S^- which can be written as a square root LL^T giving a positive semidefinite $P_3 = (W^T L)(W^T L)^T$.

We show P_1 to be positive semidefinite using the proof structure of [1] which is as follows:

$$\boldsymbol{Q}_{ab}(\boldsymbol{h}^k) = \boldsymbol{h}_a^k (\boldsymbol{P}_1)_{ab} \boldsymbol{h}_b^k$$
 (28)

 $\nu^{T} \boldsymbol{Q} \nu = \sum_{ab} \nu_{a} \boldsymbol{Q}_{ab} \nu_{b}$ $= \sum_{ab} (\boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{W})_{ab} \boldsymbol{h}_{a}^{k} \boldsymbol{h}_{b}^{k} [\frac{1}{2} \nu_{a}^{2} + \frac{1}{2} \nu_{b}^{2} - \nu_{a} \nu_{b}]$ $= \frac{1}{2} \sum_{ab} (\boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{W})_{ab} \boldsymbol{h}_{a}^{k} \boldsymbol{h}_{b}^{k} (\nu_{a} - \nu_{b})^{2}$ ≥ 0 (29)

Setting the gradient of G to zero, we obtain an expression for the minimum of G.

$$\varDelta = \frac{h^k}{(W^T S + W h^k + W^T S - x)} \odot (W^T S (W h^k - x))$$
(30)

$$\boldsymbol{h}^{k+1} = \boldsymbol{h}^{k} - \boldsymbol{\varDelta} \\ = \boldsymbol{h}^{k} \odot \frac{\boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{x} + \boldsymbol{W}^{T} \boldsymbol{S}^{-} \boldsymbol{W} \boldsymbol{h}^{k}}{\boldsymbol{W}^{T} \boldsymbol{S}^{-} \boldsymbol{x} + \boldsymbol{W}^{T} \boldsymbol{S}^{+} \boldsymbol{W} \boldsymbol{h}}$$
(31)

This is the update rule for h and similarly we can derive the update rule for w.

Projected gradient

To implement a projected gradient solver for glsNMF we require a modification of the original NMF algorithm [12] that accounts for including the covariance into the Gaussian likelihood objective. Individual objectives for W and H written in the vector form $vec(\cdot)$ (columns of a matrix stacked in a vector) are:

$$E_{W} = vec(W^{T}) \left(HH^{T} \otimes S \right) vec(W) + W's \text{ linear terms} \quad (32)$$

$$E_{H} = vec(H)^{T} (I \otimes W^{T} S W) vec(H) + H' s \text{ linear terms}, (33)$$

where \otimes is the Kronecker product defined as

$$\boldsymbol{A} \otimes \boldsymbol{B} = \begin{bmatrix} a_{11}\boldsymbol{B} & \cdots & a_{1n}\boldsymbol{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\boldsymbol{B} & \cdots & a_{mn}\boldsymbol{B} \end{bmatrix}$$
(34)

Hessian matrices in both cases tend to be well conditioned when number of training examples is larger than number of features and dimensionality of the data, and *S* is well conditioned too. The first is a usual case in the application domain of NMF. The second requires a well conditioned sample covariance, which can be provided by estimation from a sufficient population of noise or adequately modeled.

Equation (33) is structurally equivalent to the one presented by Lin [12]. The original routine for solving the subproblem of taking a step in the direction opposite to the maximal gradient $\nabla_H f(W, H)$ can be used with a single change: $W^T SW$ needs to be precomputed instead of the plain $W^T W$.

A different situation arises in the case of (32). The subroutine of H cannot be used directly and a straightforward solution is to have individual subroutines for H and W. The sufficient decrease test [12, equation 4.3] becomes:

$$(1-\sigma)\langle \nabla_{W}f(W,H), dW \rangle + \frac{1}{2}\langle dW, SdWHH^{T} \rangle \leq 0,$$

where $\langle \cdot, \cdot \rangle$ is the sum of element-wise product of two matrices.

The shape of the second term in the sum is derived from the full tensor product of the Hessian with vec(dW) taking advantage of the following identity [16]:

$$vec(\boldsymbol{A}\boldsymbol{X}\boldsymbol{B}) = (\boldsymbol{B}^T \otimes \boldsymbol{A})vec(\boldsymbol{X})$$
 (35)

The additional complexity with respect to the original projected gradient algorithm for NMF is the matrix multiplication by S. In practice we have found this not be an issue and projected gradient was considerably faster than multiplicative updates, which aligns well with what has been previously observed for NMF.



Fig. 1.

Randomly selected samples from the swimmer dataset of [6], which consists of 256 total images with all possible combination of limbs positions.





Random samples from the swimmer dataset contaminated by correlated noise. Note the salient spatial correlation with the shape of the swimmer's torso to the left of the swimmer.



Fig. 3.

Noise covariance for all 1024 image pixels shows the expected spatial correlations. The covariance is very close to identity, only the 90 by 90 close up on the right shows the correlated part. Such a covariance matrix is favorable to the conventional least squares analysis because it satisfies the assumptions of the method.



Fig. 4.

Correlated noise samples that were added to the swimmer dataset. Note the salient spatial correlation.



Fig. 5.

20 features (columns of the W matrix) as extracted by the regular NMF algorithm with the LS objective function from the dataset contaminated by correlated noise.



Fig. 6.

Features extracted by NMF from the noiseless swimmer dataset. Only 16 unique features are shown.





Fig. 7. 20 features extracted from the noisy dataset by glsNMF.