

# Positive Matrix Factorization of Large Aerosol Mass Spectrometry Datasets Using Error-Weighted Randomized Hierarchical Alternating Least Squares

Benjamin Sapper<sup>1</sup>, Daven Henze<sup>1</sup>, Manjula Canagaratna<sup>2</sup>, Harald Stark<sup>3,4</sup>, and Jose Jimenez<sup>4</sup>

<sup>1</sup>University of Colorado Boulder, 11 Engineering Dr, Boulder, CO 80309, United States

<sup>2</sup>Aerodyne Research, 45 Manning Road, Billerica, MA 01821, United States

<sup>3</sup>Center for Aerosol and Cloud Chemistry, Aerodyne Research Inc., 45 Manning Road, Billerica, MA 01821

<sup>4</sup>Department of Chemistry and Cooperative Institute for Research in Environmental Sciences (CIRES), University of Colorado, Boulder, Colorado 80309, United States

**Correspondence:** Benjamin Sapper (bsapper77@gmail.com)

**Abstract.** Weighted positive matrix factorization (PMF) has been used by scientists to find small sets of underlying factors in environmental data. However, as the size of the data has grown, increasing computational costs have made it impractical to use traditional methods for this factorization. In this paper, we present a new weighting method to dramatically decrease computational costs for these traditional algorithms. We then apply this weighting method with the Randomized Hierarchical Alternating Least Squares (RHALS) algorithm to a large environmental dataset, where we show that interpretable factors can be reproduced using these methods. We show this algorithm results in a computational speedup of 38, 67, and 634 compared to the Multiplicative Update (MU), deterministic Hierarchical Alternating Least Squares (HALS), and non-negative Alternating Least Squares (ALS) algorithms, respectively. We also investigate rotational ambiguity in the solution, and present a simple “pulling” method to rotate a set of factors. This method is shown to find alternative solutions, and in some cases, lower the weighted residual error of the algorithm.

## 1 Introduction

### 1.1 Problem Statement

Low rank matrix factorization has been widely used in data science to explain underlying factors in large datasets (~~Xie et al. (1998)~~, ~~Kim and Hopke (2007)~~, ~~Wei et al. (2016)~~) ([Xie et al., 1998](#); [Kim and Hopke, 2007](#); [Wei et al., 2016](#)). The process considers a data matrix,  $\mathbf{A}$ , of size  $m \times n$ , that is decomposed into two smaller matrices,  $\mathbf{W}$  of size  $m \times k$  and  $\mathbf{H}$  of size  $k \times n$ , where  $k \ll m, n$  and  $\mathbf{A} \approx \mathbf{WH}$ . Traditionally, principal component analysis (PCA) and singular value decomposition (SVD) have been used to find these factors (~~Kumar (2017)~~, ~~Wei et al. (2016)~~) ([Kumar, 2017](#); [Wei et al., 2016](#)). PCA finds the eigenvectors of the covariance matrix  $\mathbf{A}^T \mathbf{A}$ , which are called the principal components, representing the directions of maximum variance in the data. The vectors are ordered by how much variance they explain, and only the most important vectors are kept, which are identified as the underlying factors in the data. Closely related, the SVD finds the factorization  $\mathbf{A} = \mathbf{USV}^T$ , where  $\mathbf{U}$  and

$\mathbf{V}$  contain the left and right singular vectors, respectively, and  $\mathbf{S}$  is a diagonal matrix containing the singular values of  $\mathbf{A}$  in decreasing order. Thus, to find a low rank approximation of  $\mathbf{A}$ , one could keep only the  $k$  most significant singular values and vectors to form the truncated SVD,  $\mathbf{U}_k \mathbf{S}_k \mathbf{V}_k^T$ . Mathematically, the truncated SVD is the most optimal rank  $k$  factorization of  $\mathbf{A}$  for minimizing squared error (Eckart and Young (1936)) (Eckart and Young, 1936). However, the SVD is not optimal for some factorization problems in several ways (Paatero and Tapper (1994)) (Paatero and Tapper, 1994):

1. The SVD produces factors with negative values: for some factor analysis problems, such as finding chemical sources to air pollution data, SVD results are hard to interpret because chemical concentrations can only be positive.
2. The SVD produces orthogonal factors: most factor analysis problems do not carry the underlying constraint that the factors have to be orthogonal to each other.
3. The SVD is not fit to solve the following weighted least squares problem: suppose that accompanying the dataset  $\mathbf{A}$  is an equally sized matrix  $\Sigma$  with  $\Sigma_{ij} = \sigma_{ij}$  representing the uncertainty of the measurement for  $\mathbf{A}_{ij}$ . If  $rank(\Sigma) > 1$ , the SVD can't be scaled to find a solution minimizing the weighted residual error, which is defined as  $\sum_i \sum_j \frac{(\mathbf{A}_{ij} - \sum_k \mathbf{W}_{ik} \mathbf{H}_{kj})^2}{\sigma_{ij}^2}$  (Paatero and Tapper (1994)) (Paatero and Tapper, 1994).

Thus, positive matrix factorization (PMF) was introduced by Paatero to address these concerns (Paatero and Tapper (1993)) (Paatero and Tapper, 1994). Weighted PMF attempts to find two factor matrices  $\mathbf{W}$  and  $\mathbf{H}$  by minimizing the equation

$$\min_{\mathbf{W}, \mathbf{H}} \left\| \frac{\mathbf{A} - \mathbf{WH}}{\Sigma} \right\|_F^2 \text{ with } \mathbf{W} \geq 0, \mathbf{H} \geq 0 \quad (1)$$

In Eq. (1),  $\frac{\mathbf{A} - \mathbf{WH}}{\Sigma}$  represents elementwise division, the norm  $\|\cdot\|_F$  is the Frobenius norm, and all elements of  $\mathbf{W}$  and  $\mathbf{H}$  are constrained to be positive.

Traditional factor analysis methods are known to be computationally expensive. Steps to speed up factor analysis have thus been explored, such as randomization and the use of GPUs (Halko et al. (2011), Tan and Cao) (Halko et al., 2011; Tan et al., 2018). Recently, the authors in (Erichson et al. (2018)) Erichson et al. (2018) applied randomization to PMF, introducing a new method, Randomized Hierarchical Alternating Least Squares (RHALS) to solving the unweighted PMF problem. In this paper, we test the application of RHALS to atmospheric concentration data that contain uncertainties. Accounting for these uncertainties as regression weights, we introduce a method of externally weighting and unweighting the data, which to our knowledge is novel in its application to RHALS. We consider the accuracy and the reduced computational costs compared to other PMF algorithms commonly used in the field of atmospheric science.

## 1.2 Background

### 1.2.1 PMF2 and Paatero

The first widely accepted algorithm for PMF was derived by Paatero in Paatero (1997) using the Gauss-Newton Method (Paatero (1997)). This algorithm, called PMF2, is very commonly used with environmental data (Kim and Hopke (2007), Ulbrich et al. (2009), Massoli et al. (2018)). Paatero (Kim and Hopke, 2007; Ulbrich et al., 2009; Massoli et al., 2018). Paatero (1997) defines an

enhanced objective function, and attempts to find factor matrices  $\mathbf{W}$  and  $\mathbf{H}$  that minimize the cost function  $Q$  in

$$Q = \sum_i \sum_j \frac{(\mathbf{A}_{ij} - \sum_k \mathbf{W}_{ik} \mathbf{H}_{kj})^2}{\sigma_{ij}^2} - \alpha \sum_i \sum_k \log(\mathbf{W}_{ik}) - \beta \sum_k \sum_j \log(\mathbf{H}_{kj}) + \gamma \sum_i \sum_k \mathbf{W}_{ik}^2 + \delta \sum_k \sum_j \mathbf{H}_{kj}^2 \quad (2)$$

In Eq. (2), logarithmic penalty terms are added to penalize factor values that become too close to zero (and therefore potentially negative), and L2 regularization is added to smooth out the factors and avoid overfitting (Paatero (1997))(Paatero, 1997).  $\alpha$  and  $\beta$  control the strength of the penalty terms, while  $\gamma$  and  $\delta$  control the strength of L2 regularization. While Paatero hasn't released the source code detailing the exact algorithm to solve Eq. (2), a step-by-step guide in pseudocode as to how to use the Gauss-Newton method to solve Eq. (2) is detailed in (Lu and Wu (2004))Lu and Wu (2004).

### 1.2.2 Multiplicative Update

An alternative method for PMF was developed in (Lee and Seung (1999))Lee and Seung (1999). The multiplicative update (MU) method utilizes a special case of gradient descent where the learning rates are chosen to avoid subtraction in the gradient (Burred)(Burred, 2017). A multiplicative update estimate of a parameter  $\theta$  (either  $\mathbf{W}$  or  $\mathbf{H}$ ) is found by updates of the form (Burred)(Burred, 2017; Gillis, 2020)

$$\theta = \theta \odot \frac{\nabla_{\theta}^{-} Q(\theta)}{\nabla_{\theta}^{+} Q(\theta)} (\nabla_{\theta}^{-} Q(\theta) \oslash \nabla_{\theta}^{+} Q(\theta)) \quad (3)$$

where  $Q(\theta)$  is a cost function to be minimized,  $\nabla_{\theta}^{-}$  consists of the negative terms of the gradient of the cost function,  $\nabla_{\theta}^{+}$  consists of the positive terms of the gradient, and  $\odot$  is an element-wise multiplication and  $\oslash$  are elementwise multiplications and divisions, respectively. Since the data matrix  $\mathbf{A}$ , the uncertainties, and the factor matrices are all positive at each step of this algorithm, the factor matrices in the subsequent step will also be positive because Eq. (3) only deals with the multiplication and division of positive numbers.

It is possible to perform NMF using other forms of gradient descent - for example, the projected gradient method (PGM) sets the step size to the inverse of the maximum eigenvalue of Hessian of the cost function, and may lead to faster convergence than MU (Gillis, 2020). However, we choose to only test MU, due to its widespread use and flexibility (Gillis, 2020).

### 1.2.3 Alternating Least Squares and Hierarchical ALS

Alternating Least Squares (ALS) methods solve for the factor matrices  $\mathbf{W}$  and  $\mathbf{H}$  by iteratively updating each matrix until convergence is reached (Tan and Cao, Takács and Tikk (2012))(Tan et al., 2018; Takács and Tikk, 2012). The cost function  $Q$ , containing  $\|\frac{\mathbf{A}-\mathbf{WH}}{\Sigma}\|_F^2 \|(\mathbf{A}-\mathbf{WH}) \oslash \Sigma\|_F^2$ , is minimized by setting the partial derivatives  $\frac{\partial Q}{\partial \mathbf{W}}$  and  $\frac{\partial Q}{\partial \mathbf{H}}$  to zero, and solving for  $\mathbf{W}$  and  $\mathbf{H}$ . To satisfy the positivity constraint, negative elements in the factors are set to zero.

Nonnegative ALS has no theoretical convergence guarantee and, in some problems, may fail to converge to a feasible solution (Gillis, 2020). For this reason, Alternating Nonnegative Least Squares (ANLS) and Alternating Direction Method of Multipliers (ADMM) are interesting alternatives. In ANLS, indices of an "active set" are set to zero, and the rest are updated via

an unconstrained optimization (Kim and Park, 2011). The active set is then updated to contain the indices with the new negative factor elements. In ADMM, an auxiliary factor matrix  $\mathbf{Y}$  is formed, and an additional term is added to the cost function which penalizes the distance between the target factor matrix  $\mathbf{W}$  or  $\mathbf{H}$  and  $\mathbf{Y}$  (Gillis, 2020). Both of these methods may lead to faster and better convergence than nonnegative ALS, and are preferred by Gillis (2020). However, we find that the simple nonnegative ALS almost always converges to feasible solutions for our dataset, and we do not explore these alternative methods.

In recent years, Hierarchical Alternating Least Squares (HALS) has become increasingly popular as an efficient method for PMF (Cichocki and Phan (2009))(Cichocki and Phan, 2009). Instead of minimizing with respect to the entire factor matrices  $\mathbf{W}$  and  $\mathbf{H}$ , HALS minimizes the cost function with respect to one block, or an outer product, of individual factors at a time. The main component of the cost function is redefined as  $Q^j = \|\mathbf{R}^j - \mathbf{W}(:,j)\mathbf{H}(j,:)\|_F^2$ , where  $\mathbf{W}(:,j)$  and  $\mathbf{H}(j,:)$  are the  $j$ th factors.  $Q^j$  can be minimized for each factor  $j$  by setting the partial derivatives  $\frac{\partial Q^j}{\partial \mathbf{W}(:,j)}$  and  $\frac{\partial Q^j}{\partial \mathbf{H}(j,:)}$  to zero and solving for  $\mathbf{W}(:,j)$ , and  $\mathbf{H}(j,:)$ .

The derivation of the ALS update rules is detailed in section A of the Appendix, while the derivation of HALS is detailed in the Methods section. Section 2.

### 1.3 Reducing Computational Costs Random Projections

#### 1.3.1 Randomization

To reduce the computational costs of a matrix factorization algorithm for large datasets, randomization methods have been used as a dimension reduction technique (Erichson et al. (2018), Halko et al. (2011), Kaloorazi and Chen (2019))(Erichson et al., 2018; Halko et al., 2011). Below, we present a brief overview of the theory and results laid out in (Halko et al. (2011)) laid out in Halko et al. (2011).

When performing randomization techniques, we hope that much of the relevant information about the column space of the data matrix  $\mathbf{A}$  can be stored in just a few vectors that we can sample. This is only true if the "effective rank" of  $\mathbf{A}$  is low ( $\mathbf{A}$  only has a few nonnegligible singular values), but that is generally assumed to be the case in any PMF problem (Erichson et al. (2018))(Erichson et al., 2018). Mathematically, we seek the approximation

$$\mathbf{A} \approx \mathbf{Q}\mathbf{Q}^T\mathbf{P}\mathbf{P}^T\mathbf{A} \quad (4)$$

where the relatively few number of columns in  $\mathbf{Q}$  the matrix  $\mathbf{P}$  are orthonormal and form an approximate basis of  $\mathbf{A}$ . Choosing the columns of  $\mathbf{Q}$  to be the left singular vectors of  $\mathbf{A}$  would minimize L2 error: if the first  $k$  singular vectors were chosen, the error term  $\|\mathbf{A} - \mathbf{Q}\mathbf{Q}^T\mathbf{A}\|_2 = \sigma_{k+1}$ , with  $\sigma_{k+1}$  being the  $(k+1)^{st}$  largest singular value of  $\mathbf{A}$  (Halko et al. (2011))(Halko et al., 2011). However, it turns out that a suitable basis can also be found from random samples of the column space of  $\mathbf{A}$ .

How many samples should one take from the column space of  $\mathbf{A}$ ? Since we assume there are  $k$  underlying factors within  $\mathbf{A}$ , we essentially say that the effective rank of  $\mathbf{A}$  is  $k$  and sample  $k$  random samples from the column space of  $\mathbf{A}$ . However, with underlying uncertainties, we can write  $\mathbf{A} = \mathbf{B} + \mathbf{E}$ , where  $\mathbf{B}$  is the rank  $k$  matrix spanned by the factors for which we wish to find a basis, and  $\mathbf{E}$  is a perturbation matrix filled with the noise in  $\mathbf{A}$  (Halko et al. (2011))(Halko et al., 2011). Suppose we were

to sample from the column space of  $\mathbf{A}$  – that is, form the vector  $\mathbf{y} = \mathbf{A}\boldsymbol{\omega} = \mathbf{B}\boldsymbol{\omega} + \mathbf{E}\boldsymbol{\omega}$ . Each vector  $\mathbf{y}$  is slightly pushed out of the column space of  $\mathbf{B}$  by the term  $\mathbf{E}\boldsymbol{\omega}$ . Thus, to increase the likelihood of spanning the column space of  $\mathbf{B}$ ,  $p$  more vectors are sampled from  $\mathbf{A}$ . In practice, choosing  $p$  to be 10 or 20 is sufficient.

To construct this low rank approximation,  $k + p$  random normal samples, stored as columns of the matrix  $\boldsymbol{\Omega}$  (dimensions  $n \times (k + p)$ ) of the column space of  $\mathbf{A}$  are taken and stored in  $\mathbf{Y}$ :

$$\mathbf{Y}_{m \times (k+p)} = \mathbf{A}_{m \times n} \cdot \mathbf{\Omega}_{n \times (k+p)} \quad (5)$$

Next, the columns of  $\mathbf{Y}$  are orthonormalized using a QR decomposition to form our projection matrix  $\mathbf{Q}\mathbf{P}$ . The algorithm can now be run on the lower dimensional matrix  $\mathbf{B} = \mathbf{Q}^T \mathbf{A} \mathbf{B} = \mathbf{P}^T \mathbf{A}$ .

### 1.3.1 Parallelization

Many new algorithms use multiple cores within a computer, including Graphics Processing Units (GPUs), for computational efficiency (Tan and Cao, Nasre et al. (2013)). Tasks that are parallelizable, such as matrix multiplication, can be partitioned and run simultaneously on these cores, decreasing computing time by a factor of the number of cores used. While this paper does not explore the advantages of parallelization on the RHALS algorithm, we discuss the parts of the algorithm that can be parallelized, in particular the randomization step and the external weighting step derived in Section 2.3.

### 1.4 Non-Uniqueness of Solutions

Unlike the SVD, there is no uniqueness constraint on the factor matrices  $\mathbf{W}$  and  $\mathbf{H}$  for PMF. That is, the factorization  $\mathbf{A} = \mathbf{W}\mathbf{H}$  can also be expressed as  $\mathbf{A} = \mathbf{W}\mathbf{T}\mathbf{T}^{-1}\mathbf{H}$ , where  $\mathbf{T}$  is a "rotational" matrix and  $\hat{\mathbf{W}} = \mathbf{W}\mathbf{T}$  and  $\hat{\mathbf{H}} = \mathbf{T}^{-1}\mathbf{H}$  are the new rotated factors. Thus, to We note that  $\mathbf{T}$  does not necessarily represent a true rotation in a mathematical form, which would require  $\mathbf{T}$  to be orthogonal. To span the space of feasible solutions, previous approaches such as PMF2 have aimed to find new solutions  $\hat{\mathbf{W}}$  and  $\hat{\mathbf{H}}$  by varying  $\mathbf{T}$  (Paatero (1997)) (Paatero, 1997), or varying initializations (Ulbrich et al. (2009)) (Ulbrich et al., 2009).

The rotational matrix  $\mathbf{T}$  is a  $k \times k$  matrix where  $t_{ii}$ , a diagonal element of  $\mathbf{T}$ , represents a scaling of the  $i^{th}$  factor, and  $t_{ij}$  represents a rotation of the  $j^{th}$  factor towards the  $i^{th}$  factor in  $\mathbf{W}$ , and a rotation of the  $i^{th}$  factor away from the  $j^{th}$  factor in  $\mathbf{H}$ . For example, consider the elementary rotation matrices (Paatero et al., 2002)

$$\mathbf{T}_E = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & r & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}, \mathbf{T}_E^{-1} = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & -r & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix} \quad (6)$$

All factors remain the same, except  $\hat{\mathbf{W}}(:, 3) = \mathbf{W}(:, 3) + r\mathbf{W}(:, 2)$ , and  $\hat{\mathbf{H}}(2, :) = \mathbf{H}(2, :) - r\mathbf{H}(3, :)$ .

Regardless of whether  $r$  is positive or negative, values in either  $\mathbf{W}$  or  $\mathbf{H}$  will be pulled towards negative values. Thus if a large proportion of the factors are filled with zeros, there may be little to no pure rotations, or rotations that do not change the residual  $\|\mathbf{A} - \mathbf{WH}\|_F^2$ . Algorithms such as those developed by Paatero, where a logarithmic penalty term is added to push factor values to be more positive, will have few zeros in the factor matrices, and thus will have more rotational ambiguity (Paatero et al. (2005)) (Paatero et al., 2005). However, RHALS enforces non-negativity merely by setting negative elements to zero, and thus many values in the factor matrices may end up being zero. Thus for RHALS, there will almost certainly not exist a perfect rotational matrix  $\mathbf{T}$ , and only “approximate” rotations can be studied, in which the rotation will alter the value of the weighted error.

It is not feasible to span all possible variants that  $\mathbf{T}$  can take. Thus, the problem is often simplified to considering only positive rotations (values of  $\mathbf{T}$  greater than zero) and negative rotations (values of  $\mathbf{T}$  less than zero). A rotational program in PMF2 called FPEAK uses the parameter  $\phi$  to denote the rotation strength, with positive values leading to positive rotations in  $\mathbf{W}$  (Paatero (1997)) (Paatero, 1997). Paatero further improved this method in the Multilinear Engine (ME) algorithm, where the strength of rotation is allowed to vary between factors (Paatero and Hopke (2009)) (Paatero and Hopke, 2009). The pulling algorithm presented in (Paatero and Hopke (2009)) Paatero and Hopke (2009) is a sophisticated rotational method; more rudimentary pulling methods that mimic varying the regularization of the factor matrices are presented in (Paatero (1997)) and (Paatero et al. (2002)) Paatero (1997) and Paatero et al. (2002).

## 1.5 Scaling with Uncertainties

Recall that to account for inaccuracies in real data, we measure the ~~error~~ squared error  $\underline{Q}$  of the algorithm by dividing the residual by the standard deviation of the uncertainty of each measurement:

$$\underline{Q} = \sum_{i=1}^m \sum_{j=1}^n \left( \frac{\mathbf{A}_{ij} - \sum_{d=1}^k \mathbf{W}_{id} \mathbf{H}_{dj}}{\sigma_{ij}} \right)^2 \quad (7)$$

To account for these uncertainties, one could incorporate them into each update rule of the factor matrices, as is done in PMF2 (Paatero (1997)) (Paatero, 1997). In this paper, we refer to this as “internal weighting”. However, doing this is very computationally expensive due to ~~the size of the~~ elementwise operations with the uncertainty matrix  $\Sigma$ . ~~Furthermore, how would the uncertainties be included after the dimension of the data is reduced by a randomization step? One way to address the weighting problem introduced by dimension reduction is through an expectation maximization step layed out in (Yahaya et al. (2019)). At each iteration, the data matrix  $\mathbf{A}$  is scaled based on the uncertainties and current factors, and then compressed again into a lower dimension. However, this method is still computationally expensive—in fact, the computationally expensive step of expectation maximization eliminates any of the time benefits gained by randomization~~ Elementwise operations of large arrays are inefficient processes compared to other operations of the same computational complexity, such as matrix-vector multiplication, due to the large allocation of memory towards intermediary results (Jia et al., 2020). We thus introduce an alternate approach to weighted PMF where the data is prescaled by the uncertainties matrix, the unweighted algorithm is applied to the scaled data, and the converged factors are scaled by the uncertainties. This approach, which we will refer to as “external weighting”, thus dra-

matically reduces computational costs and allows for dimensionality reduction, since weights are not included in the update rules.

It is noted in ~~(Paatero and Tapper (1993))~~ [Paatero and Tapper \(1993\)](#) that if the weighting matrix is rank one, that is if  $1/\sigma_{ij} = \mathbf{B}_i \mathbf{C}_j$  for vectors  $\mathbf{B}$ ,  $\mathbf{C}$ , then an optimal scaling can be found. By forming the diagonal matrices  $\mathbf{D}_L$ ,  $\mathbf{D}_R$  with  $[\mathbf{D}_L]_{ii} = \mathbf{B}_i$ ,  $[\mathbf{D}_R]_{jj} = \mathbf{C}_j$ , if  $\mathbf{A} = \mathbf{W}\mathbf{H} + \mathbf{E}$ , then  $\mathbf{D}_L \mathbf{A} \mathbf{D}_R = (\mathbf{D}_L \mathbf{W})(\mathbf{H} \mathbf{D}_R) + \mathbf{D}_L \mathbf{E} \mathbf{D}_R = \hat{\mathbf{W}} \hat{\mathbf{H}} + \hat{\mathbf{E}}$ . Thus, by first finding  $\hat{\mathbf{A}} = \mathbf{D}_L \mathbf{A} \mathbf{D}_R$  and then running a PMF algorithm without weights on  $\hat{\mathbf{A}}$ , one can produce estimated factor matrices  $\hat{\mathbf{W}}$  and  $\hat{\mathbf{H}}$ , where the unscaled estimates are  $\mathbf{W} = \mathbf{D}_L^{-1} \hat{\mathbf{W}}$  and  $\mathbf{H} = \hat{\mathbf{H}} \mathbf{D}_R^{-1}$ .

If the weighting matrix  $1/\Sigma$  isn't rank one – which is very likely for environmental data – and cannot be estimated as the outer product of two vectors, there is no scaling of the previous form that can be applied to the data matrix  $\mathbf{A}$  ~~(Paatero and Tapper (1993))~~ [\(Paatero and Tapper, 1993\)](#). To address this, ~~Paatero~~ [Paatero and Tapper \(1993\)](#) presented a simple algorithm to find an approximate rank 1 factorization of  $1/\Sigma$  ~~(Paatero and Tapper (1993))~~. We present a different method where the data is scaled by the full rank matrix  $\Sigma$ , and then unscaled after the algorithm is complete. This algorithm is described in detail in Section 2.3.

### 1.5.1 [Expectation Maximization](#)

The Expectation Maximization (EM) approach was first designed for matrix factorization problems associated with missing entries (Zhang et al., 2006). Specifically, if  $A^o$  is the observed data and  $A^u$  is the unknown data within  $\mathbf{A}$ , then the EM approach seeks to find factors  $\mathbf{W}$  and  $\mathbf{H}$  that satisfy the following (Zhang et al., 2006)

$$\underset{\mathbf{W}, \mathbf{H}}{\operatorname{argmax}} \mathbb{E}(\log(\mathbb{P}(A^o, A^u | \mathbf{W}\mathbf{H})) | A^o, \mathbf{W}\mathbf{H}^{(t-1)}) \quad (8)$$

where  $\mathbf{W}\mathbf{H}^{(t-1)}$  is the product of the previous estimates of  $\mathbf{W}$  and  $\mathbf{H}$ , and  $\mathbb{P}$  is a probability measure. This problem is equivalent to running a PMF algorithm on the following adjusted matrix (Zhang et al., 2006)

$$\mathbf{A}_1 = \mathbf{C} \odot \mathbf{A} + (\mathbf{1} - \mathbf{C}) \odot \mathbf{W}\mathbf{H}^{(t-1)} \quad (9)$$

with  $\mathbf{C}_{ij} = 1$  if  $\mathbf{A}_{ij}$  is known and  $\mathbf{C}_{ij} = 0$  if  $\mathbf{A}_{ij}$  is unknown, and  $\mathbf{1}$  is a matrix of ones.

Recent work has looked into expanding on this approach to continuous weights as seen in most PMF problems of aerosol data (Yahaya et al., 2019; Yahaya, 2021; Yahaya et al., 2021). To handle the continuous case, a variation of Eq. 8 is maximized (Yahaya et al., 2019)

$$\underset{\mathbf{W}, \mathbf{H}}{\operatorname{argmax}} \mathbb{E}(\log(\mathbb{P}(\mathbf{C} \odot \mathbf{A}, (\mathbf{1} - \mathbf{C}) \odot \mathbf{A}_{theo} | \mathbf{W}\mathbf{H})) | \mathbf{C} \odot \mathbf{A}, \mathbf{W}\mathbf{H}^{(t-1)}) \quad (10)$$

where  $\mathbf{C}$  is now a weight matrix containing estimates of confidence as a value between 0 and 1 in a given data point, and  $\mathbf{A}_{theo}$  is the theoretical true data. Maximizing Eq. 10 is equivalent to running any PMF algorithm on the matrix  $\mathbf{A}_1$  formed in Eq. 9.

How should one form the confidence matrix  $\mathbf{C}$  from the uncertainties matrix  $\Sigma$ ? Yahaya (2021) suggests in scaling the weights (in this case  $1/\sigma_{ij}$ ) so that the maximum value is 1. However, previous testing has primarily focused on problems with binary weights (Yahaya et al., 2019; Yahaya, 2021; Yahaya et al., 2021).

## 1.6 Determining the Number of Factors

The quality of fit parameter, defined as  $Q = \left( \sum_i \sum_j \frac{A_{ij} - \sum_k W_{ik} H_{kj}}{\sigma_{ij}} \right)^2$ , i.e. the squared weighted residual error, can be used to determine whether a given solution either overfits or underfits the data (Ulbrich et al., 2009). If measurement errors are normally distributed, then  $Q$  will follow a  $\chi^2$  distribution with  $mn - k(m+n)$  degrees of freedom. Thus, to avoid overfitting, the number of factors in a solution are chosen such that  $Q \approx mn - k(m+n)$  (Paatero et al., 2002). When additional error is present, measuring the convergence of  $Q$  or the weighted residual error with additional factors can determine if these additional factors add much information to the model.

Another method used in determining the number of factors is the lack of rotational ambiguity of a solution (Ulbrich et al., 2009). Consider a simple case where the data matrix  $\mathbf{A}$  is the product of two rank-two matrices  $\mathbf{W} = [\mathbf{a}, \mathbf{b}]$ , and  $\mathbf{H} = [\mathbf{y}, \mathbf{z}]^T$ , with  $\mathbf{a}, \mathbf{b}, \mathbf{y}$ , and  $\mathbf{z}$  being column vectors. An exact solution can also be obtained by finding  $\mathbf{W} = [\mathbf{c}, \mathbf{d}, \mathbf{b}]$ , and  $\mathbf{H} = [\mathbf{y}, \mathbf{y}, \mathbf{z}]$ , where  $\mathbf{c} + \mathbf{d} = \mathbf{a}$ , in a process known as factor splitting (Ulbrich et al., 2009). Thus, a three factor solution introduces rotational ambiguity, as any two factors  $\mathbf{c}$  and  $\mathbf{d}$  can be chosen as long as they add up to  $\mathbf{a}$ . The same analysis can be seen by analyzing solutions with four or more factors, and we can conclude that a large amount of rotational ambiguity is a potential sign of overfitting.

## 1.7 Data

In this paper, we use the data from “Ambient Measurements of Highly Oxidized Gas-Phase Molecules during the Southern Oxidant and Aerosol Study (SOAS) 2013”, measuring highly oxidized multifunctional molecules (HOMs) over a forest site in Alabama from June 22<sup>nd</sup>, 2013 to July 7<sup>th</sup> 2013 (Massoli et al., 2018). The dataset contains mass spectra concentrations of 1,059 different ions over 27,336 different time stamps. Additionally, initial uncertainties associated with each measurement are also included. Positive Matrix Factorization was applied to the data using the PMF2 algorithm described in (Paatero (1997)) (Paatero (1997)), checking solutions from two to 10 factors, where a six factor solution was obtained (Massoli et al., 2018). We note that some of the uncertainties were artificially increased for this PMF2 analysis. The authors concluded that a significant portion of the Secondary Organic Aerosol (SOA) was anthropogenic in origin (Massoli et al., 2018) the result of interactions between biogenic and anthropogenic emissions (Massoli et al., 2018).

In this paper, we use this six factor solution as a reference solution, and test whether the RHALS algorithm can recreate formulated factors as well as those found from PMF2. For reference, the PMF2 mass spectra factors and the total mass spectra concentrations over all of the data are shown in Figure 1. The time series factors, as well as the total time series concentrations, are shown in Figure 2. Both plots show total concentration amounts over the entire time series and mass spectra, respectively. We will not discuss the chemical interpretations of the data, rather just the mathematical results from running the RHALS algorithm.



## 2 Methods

Here we present the derivation of the basic weighted HALS algorithm in Section 2.1, as well as a simple rotation algorithm in Section 2.2, our new external weighting algorithm in Section 2.3, and inclusion of L1 and L2 regularization in Section 2.4.

### 2.1 HALS Algorithm

The derivation of the unweighted algorithm for HALS is detailed in ~~(Erichson et al. (2018))~~ [Erichson et al. \(2018\)](#). Here we present a similar derivation, taking into account uncertainties in the data that act as weights. [Another derivation of weighted HALS is given in Ho \(2014\)](#).

The cost function  $Q_j$  is defined by block coordinate descent, minimizing a “block”, or outer product of individual factors, of  $\mathbf{W}$  and  $\mathbf{H}$  at a time.

$$\min Q_j = \min_{\mathbf{W}(:,j), \mathbf{H}(j,:)} \|(\mathbf{R}_j - \mathbf{W}(:,j)\mathbf{H}(j,:)) \oslash \Sigma\|_F^2 \quad (11)$$

In (11),  $\mathbf{R}_j$  is the  $j$ th residual, with

$$\mathbf{R}_j = \mathbf{A} - \sum_{i \neq j}^k \mathbf{W}(:,i)\mathbf{H}(i,:) = \mathbf{A} - \mathbf{W}\mathbf{H} + \mathbf{W}(:,j)\mathbf{H}(j,:) \quad (12)$$

$\mathbf{W}(:,j)$  is the  $j$ th column, or  $j$ th factor, of  $\mathbf{W}$ ,  $\mathbf{H}(j,:)$  is the  $j$ th row/factor of  $\mathbf{H}$ , and  $k$  is the number of factors that the algorithm is aiming to find.  $\Sigma$  contains the uncertainties associated with each element in the data matrix  $\mathbf{A}$ , and the symbol ~~is~~ [element-wise](#) ~~is~~ [elementwise](#) division.

To derive update rules for HALS, partial derivatives of Eq. (11) are taken with respect to the factors  $\mathbf{W}(:,j)$  and  $\mathbf{H}(j,:)$ . With  $\Sigma$  present, this can become tricky, so we present a variation on the derivation presented in ~~(Erichson et al. (2018))~~ [Erichson et al. \(2018\)](#) by considering just a row and column of the weighted residual

$$Q_j^i = \|(\mathbf{R}_j(i,:) - \mathbf{W}(i,j)\mathbf{H}(j,:))\Sigma_i^{-1}\|_F^2 \quad (13)$$

$$Q_j^p = \|\Sigma_p^{-1}(\mathbf{R}_j(:,p) - \mathbf{H}(j,p)\mathbf{W}(:,j))\|_F^2 \quad (14)$$

In Eq. (13) and Eq. (14),  $\Sigma_i$  are  $\Sigma_p$  are diagonal matrices with the diagonal elements corresponding to the elements of the  $i$ th row and  $p$ th column of  $\Sigma$ , respectively. Expanding Eq. (13) and Eq. (14) using the fact that  $\|\mathbf{X}\|_F^2 = \text{Tr}(\mathbf{X}^T\mathbf{X})$ , where ~~takes~~ [denotes](#) the trace of the matrix, we get

$$Q_j^i = \text{Tr}(\Sigma_i^{-1}\mathbf{R}_j^T(i,:) \mathbf{R}_j(i,:) \Sigma_i^{-1} - 2\mathbf{W}(i,j)\Sigma_i^{-1}\mathbf{R}_j^T(i,:) \mathbf{H}(j,:) \Sigma_i^{-1} + \mathbf{W}(i,j)^2 \Sigma_i^{-1} \mathbf{H}^T(j,:) \mathbf{H}(j,:) \Sigma_i^{-1}) \quad (15)$$

$$Q_j^p = \text{Tr}(\mathbf{R}_j^T(:,p) \Sigma_p^{-1} \Sigma_p^{-1} \mathbf{R}_j(:,p) - 2\mathbf{H}(j,p) \mathbf{R}_j^T(:,p) \Sigma_p^{-1} \Sigma_p^{-1} \mathbf{W}(:,j) + \mathbf{H}(j,p)^2 \mathbf{W}^T(:,j) \Sigma_p^{-1} \Sigma_p^{-1} \mathbf{W}(:,j)) \quad (16)$$

~~Taking a derivative~~ Differentiating with respect to  $\mathbf{W}(i, j)$  and  $\mathbf{H}(j, p)$ ,

$$\frac{\partial Q_j^i}{\partial \mathbf{W}(i, j)} = -2Tr(\Sigma_i^{-1} \mathbf{R}_j^T(i, :) \mathbf{H}(j, :) \Sigma_i^{-1}) + 2\mathbf{W}(i, j) Tr(\Sigma_i^{-1} \mathbf{H}^T(j, :) \mathbf{H}(j, :) \Sigma_i^{-1}) \quad (17)$$

$$\frac{\partial Q_j^p}{\partial \mathbf{H}(j, p)} = -2Tr(\mathbf{R}_j^T(:, p) \Sigma_p^{-1} \Sigma_p^{-1} \mathbf{W}(:, j)) + 2\mathbf{H}(j, p) Tr(\mathbf{W}^T(:, j) \Sigma_p^{-1} \Sigma_p^{-1} \mathbf{W}(:, j)) \quad (18)$$

260 To eliminate the matrix traces, it is easy to show that Eq. (17) and Eq. (18) can be rewritten as

$$\frac{\partial Q_j^i}{\partial \mathbf{W}(i, j)} = -2\mathbf{R}_j(i, :)(\mathbf{H}^T(j, :)/\underline{\oslash}(\Sigma(\underline{i}, :)^T \odot \Sigma(\underline{i}, :)^T)) + 2\mathbf{W}(i, j)\mathbf{H}(j, :)(\mathbf{H}^T(j, :)/\underline{\oslash}(\Sigma(\underline{i}, :)^T \odot \Sigma(\underline{i}, :)^T)) \quad (19)$$

$$\frac{\partial Q_j^p}{\partial \mathbf{H}(j, p)} = -2\mathbf{R}_j^T(:, p)(\mathbf{W}(:, j)/\underline{\oslash}(\Sigma_p(:, p) \odot \Sigma_p(:, p))) + 2\mathbf{H}(j, p)\mathbf{W}^T(:, j)(\mathbf{W}(:, j)/\underline{\oslash}(\Sigma_p(:, p) \odot \Sigma_p(:, p))) \quad (20)$$

where  $\underline{\oslash}$  again represents element-wise division and  $\odot$  represents element-wise multiplication. Setting Eq. (19) and Eq. (20) to zero and solving for the factor values yields

$$265 \quad \mathbf{W}(i, j) = \frac{\mathbf{R}_j(i, :)(\mathbf{H}^T(j, :)/(\Sigma_i^T \odot \Sigma_i^T))}{\mathbf{H}(j, :)(\mathbf{H}^T(j, :)/(\Sigma_i^T \odot \Sigma_i^T))} \frac{\mathbf{R}_j(i, :)(\mathbf{H}^T(j, :) \odot (\Sigma(i, :)^T \odot \Sigma(i, :)^T))}{\mathbf{H}(j, :)(\mathbf{H}^T(j, :) \odot (\Sigma(i, :)^T \odot \Sigma(i, :)^T))} \quad (21)$$

$$\mathbf{H}(j, p) = \frac{\mathbf{R}_j^T(:, p)(\mathbf{W}(:, j)/(\Sigma_p \odot \Sigma_p))}{\mathbf{W}^T(:, j)(\mathbf{W}(:, j)/(\Sigma_p \odot \Sigma_p))} \frac{\mathbf{R}_j^T(:, p)(\mathbf{W}(:, j) \odot (\Sigma(:, p) \odot \Sigma(:, p)))}{\mathbf{W}^T(:, j)(\mathbf{W}(:, j) \odot (\Sigma(:, p) \odot \Sigma(:, p)))} \quad (22)$$

Substituting Eq. (12) into Eq. (21) and Eq. (22) yields the update rules:

$$\mathbf{W}(i, j) \leftarrow \left[ \mathbf{W}(i, j) + \frac{\mathbf{A}(i, :)(\mathbf{H}^T(j, :)/(\Sigma_i^T \odot \Sigma_i^T)) - \mathbf{W}(i, :)\mathbf{H}(\mathbf{H}^T(j, :)/(\Sigma_i^T \odot \Sigma_i^T))}{\mathbf{H}(j, :)(\mathbf{H}^T(j, :)/(\Sigma_i^T \odot \Sigma_i^T))} \frac{\mathbf{A}(i, :)(\mathbf{H}^T(j, :) \odot (\Sigma(i, :)^T \odot \Sigma(i, :)^T)) - \mathbf{W}(i, :)\mathbf{H}(j, :)(\mathbf{H}^T(j, :) \odot (\Sigma(i, :)^T \odot \Sigma(i, :)^T))}{\mathbf{H}(j, :)(\mathbf{H}^T(j, :) \odot (\Sigma(i, :)^T \odot \Sigma(i, :)^T))} \right] \quad (23)$$

$$\mathbf{H}(j, p) \leftarrow \left[ \mathbf{H}(j, p) + \frac{\mathbf{A}^T(:, p)(\mathbf{W}(:, j)/(\Sigma_p \odot \Sigma_p)) - \mathbf{H}^T(:, p)\mathbf{W}^T(\mathbf{W}(:, j)/(\Sigma_p \odot \Sigma_p))}{\mathbf{W}^T(:, j)(\mathbf{W}(:, j)/(\Sigma_p \odot \Sigma_p))} \frac{\mathbf{A}^T(:, p)(\mathbf{W}(:, j) \odot (\Sigma(:, p) \odot \Sigma(:, p))) - \mathbf{H}^T(:, p)\mathbf{W}^T(:, j)(\mathbf{W}(:, j) \odot (\Sigma(:, p) \odot \Sigma(:, p)))}{\mathbf{W}^T(:, j)(\mathbf{W}(:, j) \odot (\Sigma(:, p) \odot \Sigma(:, p)))} \right] \quad (24)$$

270 where  $[v]_+$  projects all negative values of  $v$  to 0. In practice, the authors in [Erichson et al. \(2018\)](#) [Erichson et al. \(2018\)](#) utilize the following simplified form of Eq. (23) and Eq. (24):

$$\mathbf{W}(i, j) \leftarrow \left[ \mathbf{W}(i, j) - \frac{\nabla_{\mathbf{W}_{ij}} Q_j^i}{\nabla_{\mathbf{W}_{ij}}^2 Q_j^i} \right]_+ \quad (25)$$

$$\mathbf{H}(j, p) \leftarrow \left[ \mathbf{H}(j, p) - \frac{\nabla_{\mathbf{H}_{jp}} Q_j^p}{\nabla_{\mathbf{H}_{jp}}^2 Q_j^p} \right]_+ \quad (26)$$

Thus, one can add additional auxiliary functions to the cost function  $Q$ , such as regularization and rotation terms, and add  
275 them to the update rules based on the new gradient and Hessian values. Furthermore, writing the update rules as in Eq. (25) and Eq. (26) includes the calculation of the projected gradient. This can be used as a stopping condition criterion that avoids the computational costs of calculating other convergence statistics ([Erichson et al. \(2018\)](#)) ([Erichson et al., 2018](#)).

We note that when the uncertainties are equal to 1,  $\Sigma$  can be disregarded in Eq. (23) and Eq. (24), and the update rules are identical to those in ([Erichson et al. \(2018\)](#)) [Erichson et al. \(2018\)](#).

280 In Eqs. 23 and 24, the Hessians  $\mathbf{H}(\mathbf{H}^T(j, :) \odot (\Sigma^T(i, :) \odot \Sigma^T(i, :)))$  and  $\mathbf{W}^T(\mathbf{W}(:, j) \odot (\Sigma(:, p) \odot \Sigma(:, p)))$  should be found prior to multiplication by  $\mathbf{W}(i, :)$  and  $\mathbf{H}^T(:, p)$ , respectively, to minimize computational costs. We do not preallocate the products  $\mathbf{WH}$  and  $\mathbf{H}^T \mathbf{W}^T$ , which is the same implementation as in [Erichson et al. \(2018\)](#).

## 2.2 Rotational Considerations

Suppose ~~a factor is that a solution is mostly~~ interpretable by a scientist running a PMF algorithm, but ~~a particular component~~  
285 ~~is lacking or overrepresented~~ some aspects of the solution appear unrealistic, e.g. a factor is zero during a period in which it is expected to be present, or 2 factors appear mixed in their time series and/or spectra. The scientist might be interested in a rotated solution that ~~includes the complete factor. As detailed in (Paatero and Hopke (2009))~~ adjusts this. As mentioned in Section 1.4, we do not attempt to constrain these “rotations” to be norm preserving. However, it is possible to find approximate rotations. As detailed in [Paatero and Hopke \(2009\)](#), for a specific factor value  $\mathbf{W}_{ij}$  or  $\mathbf{H}_{jp}$ , an auxiliary term can be added to  
290 the cost function to pull the component towards a set value  $\mathbf{W}_{ij}^*$  or  $\mathbf{H}_{jp}^*$ . Defining

$$Q_{aux_W}^{i,j} = s(\mathbf{W}_{ij} - \mathbf{W}_{ij}^*)^2 \text{ or } Q_{aux_H}^{p,j} = s(\mathbf{H}_{jp} - \mathbf{H}_{jp}^*)^2 \quad (27)$$

with  $s$  determining the strength of the pull, we find  $\nabla_{\mathbf{W}_{ij}} Q_{aux_W}^{i,j} = 2s(\mathbf{W}_{ij} - \mathbf{W}_{ij}^*)$ ,  $\nabla_{\mathbf{W}_{ij}}^2 Q_{aux_W}^{i,j} = 2s$ ,  $\nabla_{\mathbf{H}_{jp}} Q_{aux_H}^{p,j} = 2s(\mathbf{H}_{jp} - \mathbf{H}_{jp}^*)$ , and  $\nabla_{\mathbf{H}_{jp}}^2 Q_{aux_H}^{p,j} = 2s$ . Adding these terms to the gradient and Hessian in Eq. (25) and Eq. (26), we can derive new update rules for  $\mathbf{W}_{ij}$  and  $\mathbf{H}_{jp}$

295 To probe possible rotations of entire factors, we introduce “pulling equations” to the cost function, which pull the elements of  $\mathbf{W}$  and  $\mathbf{H}$  to the desired rotations:

$$Q_{aux_H}^{(j)} = a \left( \sum_{i=1}^n (1 - \mathbf{H}_{ji}) \right)^2 \quad (28)$$

$$Q_{aux_W}^{(j)} = b \left( \sum_{i=1}^m \mathbf{w}_{ij} \right)^2 \quad (29)$$

300 These equations correspond to  $r < 0$ , with both  $a$  and  $b$  positive. For  $r > 0$ ,  $\mathbf{W}$  and  $\mathbf{H}$  would be interchanged. These pulling equations are similar to those introduced in Paatero et al in ~~(Paatero et al. (2002))~~ [Paatero et al. \(2002\)](#), although their meanings are slightly changed. Simply, we aim to pull values of  $\mathbf{H}$  towards more positive values and values of  $\mathbf{W}$  towards 0, which mirrors the effect of setting  $r$  to be less than zero in Eq. (6).

Taking derivatives of Eq. (28) and Eq. (29),

$$305 \quad \nabla_{\mathbf{w}_{ij}} Q_{aux_W}^{j,i} = b \mathbf{w}_{ij} \quad (30)$$

$$\nabla_{\mathbf{w}_{ij}}^2 Q_{aux_W}^{j,i} = b \quad (31)$$

$$\nabla_{\mathbf{h}_{jp}} Q_{aux_H}^{j,p} = -a(1 - \mathbf{h}_{jp}) \quad (32)$$

310

$$\nabla_{\mathbf{h}_{jp}}^2 Q_{aux_H}^{j,p} = a \quad (33)$$

we add these terms to the gradient and Hessian in Eq. (25) and Eq. (26) to ~~receive-form~~ new update rules.

Optimally, we would expect these equations to pull the values of  $\mathbf{H}$  up and the values of  $\mathbf{W}$  towards zero, and to reach different, potentially better solutions.

## 315 2.3 External Weighting

To perform external weighting, we first find  ~~$\hat{\mathbf{A}} = \mathbf{A} ./ \Sigma$~~   $\hat{\mathbf{A}} = \mathbf{A} \oslash \Sigma$ , i.e., we divide the data elementwise by the uncertainties. ~~RHALS is then~~ We note that the uncertainty matrix  $\Sigma$  must only contain nonzero real entries. Thus, it can not handle problems with binary weights, such as PMF problems with missing entries. One approach to PMF with binary weights is the Expectation-Maximization (EM) approach detailed in Section 1.5.1 (Yahaya et al., 2021; Zhang et al., 2006).

320 After elementwise division, RHALS is applied to  $\hat{\mathbf{A}}$  to form the estimated scaled factor matrices  $\hat{\mathbf{W}}$  and  $\hat{\mathbf{H}}$ . The unscaled estimates  $\mathbf{W}$  and  $\mathbf{H}$  are found by the relation

$$\mathbf{WH} = (\hat{\mathbf{W}}\hat{\mathbf{H}}) \odot \Sigma \quad (34)$$

where  $\odot$  represents elementwise multiplication.  $\mathbf{W}$  and  $\mathbf{H}$  can then be found iteratively via alternating least squares:

$$\mathbf{W} = ((\hat{\mathbf{W}}\hat{\mathbf{H}}) \odot \Sigma) \mathbf{H}^T (\mathbf{H}\mathbf{H}^T)^{-1} \text{ or } \mathbf{W} = ((\hat{\mathbf{W}}\hat{\mathbf{H}}) \odot \Sigma) \mathbf{H}^\dagger \quad (35)$$

$$\mathbf{H} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T ((\hat{\mathbf{W}} \hat{\mathbf{H}}) \odot \Sigma) \text{ or } \mathbf{H} = \mathbf{W}^\dagger ((\hat{\mathbf{W}} \hat{\mathbf{H}}) \odot \Sigma) \quad (36)$$

Here,  $\mathbf{W}^\dagger$  and  $\mathbf{H}^\dagger$  denotes the pseudoinverses of  $\mathbf{W}$  and  $\mathbf{H}$ , and all negative elements are set to zero after the factor matrix is updated. Mathematically, these methods are identical, as long as the rank of the factor matrices is equal to  $k$ , and a short proof of this is detailed in section B of the Appendix. Of course, if the factor matrices were of lower rank (or contained very small singular values), then a lower rank factorization should first be found instead. Since both algorithms are identical mathematically, one could be favorable if it provided a speed advantage, depending on the computational efficiency of the pseudoinverse algorithm called. In Matlab, a Krylov subspace algorithm is used for calculating the pseudoinverses, and both update rules are  $\mathcal{O}(mnk)$  (Feng et al., 2018). Running this code in Matlab on a single CPU, we found the update rules using the pseudoinverses were ~~less-computationally-expensive~~ faster.

To begin the iteration, we have to initiate either  $\mathbf{W}$  or  $\mathbf{H}$ . The most intuitive way to do this is by assuming the weighted factors are similar to the unweighted factors, and setting  $\mathbf{H}$  (or  $\mathbf{W}$ ) to the weighted values, but scaled to the magnitude of the original data. Since each entry of the data matrix is approximated as  $\mathbf{A}_{ij} = \sum_{p=1}^k \mathbf{W}_{ip} \mathbf{H}_{pj}$ , and assuming relatively equal magnitudes for  $\mathbf{W}$  and  $\mathbf{H}$ , the appropriate scale factor should be  $\sqrt{\frac{\text{mean}(\mathbf{A})}{k}}$ . When  $\mathbf{W}$  and  $\mathbf{H}$  contain values with greatly differing magnitudes, perhaps a different scaling factor should be used, although this has not been extensively tested. Thus, the algorithm should be initialized by

$$\mathbf{H}_0 = \sqrt{\frac{\text{mean}(\mathbf{A})}{k}} \frac{\mathbf{H}}{\text{mean}(\mathbf{H})} \quad (37)$$

Measuring the change in  $W$  and  $H$  between iterations, the algorithm typically converges relatively quickly - within 20 to 40 iterations. ~~However, the algorithm may take longer for a larger~~ In practice, one can use L2 regularization in the external weighting steps, equal to 0.01 for our data, as the least squares method may become increasingly ill-posed as the number of factors and in certain cases. More work is needed regarding the different convergence issues that may arise increases. This value may need to be altered based on the magnitude of the values in the data as well as the number of factors. However, we found that adding L2 regularization lowered the similarity of the factors to the given factors from the solution using PMF2.

We used this method for all externally weighted algorithms tested in Section 4. However, theoretically, any PMF algorithm could be used for the post processing step. This may become relevant, since as noted in Section 2.2.3, the nonnegative ALS method described above may have convergence issues for certain factorization problems (Gillis, 2020).

## 2.4 Regularization

As in the RHALS algorithm presented in ~~(Erichson et al. (2018))~~ Erichson et al. (2018), we add L1 and L2 regularization. L1 regularization is added to the cost function through the L1 norm,  $|||_1$ , which is the sum of the absolute value of the components in a given row/column in the factor matrices. L2 regularization is added through the L2, or Euclidean, norm,  $|||_2^2$ . L1 regularization is typically added to reduce sparsity, while L2 regularization is added to control the Euclidean norms of the factors

as well as avoid overfitting (~~Erichson et al. (2018)~~) ill-posed problems (Erichson et al., 2018; Hansen and O’Leary, 1993). The cost function now becomes

$$\frac{\|(\mathbf{R}_j - \mathbf{W}(:,j)\mathbf{H}(j,:))\|_F^2}{\Sigma} + \alpha\|\mathbf{W}(:,j)\|_1 + \beta\|\mathbf{H}(j,:)\|_1 + \gamma\|\mathbf{W}(:,j)\|_2^2 + \delta\|\mathbf{H}(j,:)\|_2^2$$

360

$$\|(\mathbf{R}_j - \mathbf{W}(:,j)\mathbf{H}(j,:))\|_F^2 + \alpha\|\mathbf{W}(:,j)\|_1 + \beta\|\mathbf{H}(j,:)\|_1 + \gamma\|\mathbf{W}(:,j)\|_2^2 + \delta\|\mathbf{H}(j,:)\|_2^2 \quad (38)$$

In 38,  $\alpha, \beta, \gamma$ , and  $\delta$  control the amount of regularization that is added. Typically  $\alpha$  and  $\beta$ , as well as  $\gamma$  and  $\delta$ , are set equal.  $\nabla_{\mathbf{W}_{ij}} Q_j^i = \beta + 2\gamma\mathbf{W}_{ij}$ ,  $\nabla_{\mathbf{W}_{ij}}^2 Q_j^i = 2\gamma$ , and  $\nabla_{\mathbf{H}_{jp}} Q_j^p = \alpha + 2\delta\mathbf{H}_{jp}$ ,  $\nabla_{\mathbf{H}_{jp}}^2 Q_j^p = 2\delta$ . These terms are added to the update rules in Eq. (25) and Eq. (26).

365 What should  $\alpha, \beta, \gamma$ , and  $\delta$  in Eq. (38) be in practice? Optimal parameter values are often found using an L curve analysis, which measures the tradeoff between minimizing a norm and minimizing the residual error (~~Hansen and O’Leary (1993)~~) (Hansen and O’Leary, 1993).

In the SOAS dataset, the average data value is  $1.9969 \times 10^{-5}$ . When external weighting is applied, the average value of  $\frac{\mathbf{A} \odot \Sigma}{\Sigma}$  is 1.4851. L curves below are thus plotted using two different scales: one with a scale of  $10^{-5}$ , and one with a scale of 1. The deterministic, unweighted algorithm is first applied to a small dataset consisting of the first 100 rows and columns of the SOAS dataset, scaled to contain different average values, to produce the L curves shown in Figure 3 Values of the regularization parameters are listed on the graphs, with “e” representing the exponential function with base 10, and the number to the right is the exponent (i.e.  $e03=10^3$ ). The optimal regularization value is one that is located in the bend of the graph, meaning it minimizes a solution’s norm while having low residual error. In our results, we put preference on minimizing the latter, so we choose regularization values towards the left side of the bend. Analyzing the figures, the L1 regularization parameters  $\alpha$  and  $\beta$  are chosen between  $1 \times 10^{-6}$  and  $1 \times 10^{-5}$  for the scale of  $10^{-5}$ , and between 0.1 and 1 for the scale of 1. The L2 regularization parameters  $\gamma$  and  $\delta$  are chosen around .001 for the scale of  $10^{-5}$ , and 100 for the scale of 1. We note that choosing much smaller regularization values does not drastically increase the norms of the solution, suggesting that regularization is not especially necessary for this factorization. However, for more ill-posed problems, the L1 and L2 norms may become extremely large as the amount of regularization tends to zero (Hansen and O’Leary, 1993).

370

375

380

As a rule of thumb, the L1 regularization parameters should be chosen to be on the same order of magnitude as the data, and the L2 regularization parameters should be chosen to be around 100 times that. Note that neither the data matrix size nor the number of factors affect these optimal regularization parameter values. External and internal weighting will also not affect these values, as the magnitudes of the gradient and the Hessian will be around the same value.

385 An increasingly popular alternative to traditional regularization is nuclear norm regularization, which can be applied to matrix factorization without the need for rank constraints (Hu et al., 2013; Sun and Mazumder, 2013; Fornasier et al., 2011).

The nuclear norm is defined as

$$\|\mathbf{X}\|_* = \sum_{i=1}^n \sigma_i(\mathbf{A})$$

where  $\sigma_i(\mathbf{A})$  is the  $i^{th}$  singular value of  $\mathbf{A}$ . It is the convex envelope to the rank function, and thus finding a PMF solution that minimizes the nuclear norm also has a minimal rank (Hu et al., 2013). Two possible ways to apply nuclear norm regularization to weighted PMF are by the Alternating Direction Method of Multipliers (ADMM), as demonstrated in Sun and Mazumder (2013), and by reconstruction of the nuclear norm into a Frobenius norm, which can then be treated as L2 regularization (Fornasier et al., 2011). However, both these approaches involve key computations with the large low rank product  $\mathbf{WH}$ , and thus may be computationally expensive to implement. Furthermore, both algorithms still involve some arbitrary choice of rank by requiring a preset amount of nuclear norm regularization. In traditional PMF, factor profiles are optimized from a prechosen amount of factors, and thus we center our results around this approach.

### 3 Computational Details

The computer used for these calculations has an 11th Gen Intel® Core™ i7-1165G7 dual quad-core CPU with a speed of 2.80GHz. MATLAB™ R2021a is used for all calculations. While no algorithms detailed in this paper explicitly utilize multicore processing, they may rely on built in Matlab functions (such as `svd()`) that utilize multicore processing for increased efficiency.

### 4 Results

We measure the accuracy of RHALS using the weighted residual error of the final algorithm again defined as

$\sqrt{\left(\sum_i \sum_j \frac{\mathbf{A}_{ij} - \sum_k \mathbf{W}_{ik} \mathbf{H}_{kj}}{\sigma_{ij}}\right)^2} / \sqrt{\sum_i \sum_j \left(\frac{\mathbf{A}_{ij} - \sum_k \mathbf{W}_{ik} \mathbf{H}_{kj}}{\sigma_{ij}}\right)^2}$ , the correlation coefficients between the time series factors produced by RHALS and those of the reference solution, and the cosine similarity between the mass spectra profiles produced by RHALS and the reference solution.

The MU algorithm is initialized with random numbers, while the ALS and HALS algorithms are initialized by ~~a SVD through an the nonnegative double SVD (NNDSVD)~~ approach detailed in ~~(Boutsidis and Gallopoulos (2008))~~ Boutsidis and Gallopoulos (2008). Note that the latter approach will still vary between random seeds, since the non-negative SVD is also found from utilizing the randomization technique in Section 1.3.1.

All algorithms cease updates when a stopping criterion is reached, or the maximum number of iterations is hit. For the MU and ALS algorithms, the value of the cost function is calculated after each iteration, and the algorithm halts when the percent change in this value, divided by the initial cost function value, is less than a set tolerance. The HALS algorithms use a different metric: the percent change in the projected gradient, which is detailed in ~~(Lin (2007))~~ Lin (2007). The projected gradient is

415 defined as ~~(Erichson et al., (2018))~~ (Erichson et al., 2018)

$$\nabla_{\mathbf{W}_{ij}}^P \begin{cases} \nabla_{\mathbf{W}_{ij}} & \text{if } \mathbf{W}_{ij} > 0 \\ \min(0, \nabla_{\mathbf{W}_{ij}}) & \text{if } \mathbf{W}_{ij} = 0 \end{cases} . \quad (39)$$

and similarly for  $\nabla_{\mathbf{H}_{jp}}^P$ . When the percentage change in  $\sum_i \sum_j (\nabla_{\mathbf{W}_{ij}}^P)^2 + \sum_j \sum_p (\nabla_{\mathbf{H}_{jp}}^P)^2$  (divided by its initial value) is less than a specified tolerance, the algorithm halts. Below, the maximum number of iterations allowed is set to be 100 (which is rarely reached by any algorithm), and the tolerance is set to be  $1 \times 10^{-4}$  for the HALS and ALS algorithms. A tolerance of  
420  $10^{-5}$  is used for the MU algorithm to account for a larger initial cost from the random initialization.

When comparing factors from different algorithms, it is important to note that the ratio of values between a time series factor and a mass spectra factor might be slightly different. Consider two algorithms that compute the identical factorizations

$$\mathbf{W}_1 \mathbf{H}_1 = \sum_{i=1}^k \mathbf{W}_i^{(1)} \mathbf{H}_i^{(1)} \text{ and } \mathbf{W}_2 \mathbf{H}_2 = \sum_{i=1}^k \mathbf{W}_i^{(2)} \mathbf{H}_i^{(2)} = \sum_{i=1}^k \left( c_i \mathbf{W}_i^{(1)} \right) \left( \frac{1}{c_i} \mathbf{H}_i^{(1)} \right)$$

with  $\mathbf{W}_i$  being a column of  $\mathbf{W}$ ,  $\mathbf{H}_i$  being a row of  $\mathbf{H}$ , and  $c_k$  being an arbitrary positive constant. If we were to plot and  
425 compare  $\mathbf{W}_i^{(1)}$  against  $\mathbf{W}_i^{(2)}$ , we would see that  $\mathbf{W}_i^{(2)}$  is only  $c_i$  multiplied by the magnitude of  $\mathbf{W}_i^{(1)}$ , concluding one algorithm predicts a stronger signal for factor  $i$ , despite the fact that the factorizations are identical. Thus, when comparing individual factors between algorithms, the time series concentrations are cumulative over the entire mass spectra, and the mass spectra concentrations are cumulative over the entire time series (as also presented with the PMF2 results). Specifically, we plot the time series of factor  $i$  as the sum over the columns of the outer product  $\mathbf{W}(:, i) \mathbf{H}(i, :)$ , and we plot the mass spectra of factor  $i$   
430 as the sum over the rows.

To avoid this issue, one can add scaling coefficients to rescale the factors at each iteration so that they will approximately have the same magnitude ~~(Lu and Wu (2004))~~ (Lu and Wu, 2004). This has not yet implemented in our code.

#### 4.1 Computational Efficiency

~~First, to show-~~

435 Table 1 shows the number of operations required at each step (update of  $\mathbf{W}$  and  $\mathbf{H}$ ). Preprocessing steps such as finding  $\mathbf{A} \oslash \mathbf{\Sigma}$  and  $\mathbf{\Sigma} \odot \mathbf{\Sigma}$  are excluded from the table. External Weighting eliminates almost all elementwise operations needed, which is a slow memory bound operation (Jia et al., 2020). However, systems with an abundance of free memory and/or GPUs may find that internal weighting methods are sufficiently quick for small or medium size problems. Furthermore, the performance of matrix operations within these algorithms may vary between programming languages and libraries, such as between the  
440 numpy, numexpr, and Theano libraries in Python (Bergstra et al., 2010).

Analyzing Table 1 allows us to see why, in practice, it takes longer for internally weighted ALS and HALS to converge than internally weighted MU, as the former will have more matrix multiplication operations and elementwise operations for our target rank,  $k = 6$ . However, when the algorithms are externally weighted, ALS and HALS become much more computationally efficient, allowing runtime to diminish as well.



445 [To show numerically](#) how the cost of RHALS scales with the size of the system, we consider the size of the data matrix vs the run time of the RHALS algorithm, as well as the deterministic internally weighted HALS and MU algorithms. Figure 4 shows how the computational costs of the different algorithms vary with the number of rows (and columns) in a square data matrix. The number of factors is set to five, and the data is formed by multiplying together two rank five matrices filled with sampled values from the PMF2 solution of the SOAS dataset. Uncertainty values are also randomly sampled from the given  
 450 uncertainty data, and noise is added to the data matrix through random normal values that have a standard deviation of the uncertainties. L1 regularization is set to 1 for all factors, and L2 regularization is set to 50.

As one can see by looking at Figure 4, the logarithms of size and run time look to be [proportional](#)[linear](#), and thus data size and run time look to be proportional with all three algorithms. RHALS clearly outperforms MU and HALS in terms of efficiency, taking mere seconds to run for each matrix size, about 1/20 of the time of the MU algorithm and 1/100 of the time of  
 455 the deterministic HALS algorithm. It should also be noted that the randomization step, which accounts for about five percent of the algorithm computational costs, can be parallelized ([Erichson et al. \(2018\)](#)) [using Graphics Processing Units \(GPUs\)](#) ([Erichson et al., 2018](#)). Additionally, since the post algorithm external weighting step involves just matrix multiplication and an SVD, it can also be parallelized ([Lahabar and Narayanan \(2009\)](#)) ([Lahabar and Narayanan, 2009](#)). This post weighting step accounts for around 60% of the algorithm’s computational costs.

## 460 4.2 Simple Case

We now present the RHALS algorithm over a small test case. To do this, we will form a data matrix by combining underlying “true” factors and random noise. Thus, first we form a 5,952 by 400 matrix using the first three factors from the PMF2 solution to the SOAS dataset. We choose the time stamps between 9385 and 15336 from the PMF2 time series factor matrix  $\mathbf{W}_{\text{PMF2}}$ , since the PMF2 factors are clearly distinguishable from each other on this time interval. We choose the first 400 mass spectra  
 465 profiles from the mass spectra factor matrix  $\mathbf{H}_{\text{PMF2}}$ , since the bulk of the mass spectra concentrations lie within this range.

A data matrix is then formed by combining the factors and adding random normal noise with a standard deviation equal to that of the uncertainty of the data:

$$\mathbf{A}_{\text{test}}(i, j) = \max(0, \mathbf{W}_{\text{test}}(i, :) \mathbf{H}_{\text{test}}(:, j) + N(0, \Sigma(i, j))) \quad (40)$$

Next, we run RHALS over 20 trials, each with a different initialization, in order to obtain many different solutions that exist  
 470 in the solution space. Here, initializations are formed using the randomized SVD method described earlier, which works well enough for generating different solutions. Figure 5 shows the weighted residual error of the RHALS algorithm over 20 trials.

The average weighted error over 20 trials is  $1.4749 \times 10^3$ , and the algorithm took an average time of 0.0561 seconds to converge with an average of 48.15 block coordinate descent steps per trial. As one can see, solutions can vary over different trials, and a smaller weighted error could be used as justification of one solution over another.

475 Next, we compare the similarities of the converged factors and the original factors that formed the test case. The averages of the similarities between the three factors are shown in Figure 6. A cosine similarity is used for the mass spectra, while the correlation coefficient is used for the time series.

As one can see, the RHALS algorithm recreates the mass spectra and time series factors almost perfectly for all trials tested, regardless of the weighted residual error. In practice, a cosine similarity over .95 or a correlation coefficient over .90 between factor profiles represent almost identically interpretable solutions. Any trial is viable to be chosen as a “good” solution. Each solution yields a mass spectra similarity over 0.994, and a time series correlation over 0.974. There is low variance in the similarity metrics among trials, and higher weighted error thus corresponds with less important parts of the factor.

### 4.3 Large Dataset

#### 4.3.1 Comparing Different Algorithms

Next, we analyze the complete dataset, and compare the RHALS factors to all PMF2 factors. Table 2 shows a table of diagnostics for the HALS, RHALS, ALS, and MU algorithms applied to the complete SOAS dataset, averaged over three trials. As one can see, algorithms with external weighting demonstrate a dramatic reduction in computational costs, albeit at a larger weighted residual error. The ALS algorithms were by far the slowest, and the internally weighted HALS algorithm yielded the lowest average weighted residual error. The internally weighted MU algorithm produced comparable results to the internally weighted HALS algorithm in less than two thirds of the time; however, both the internally and externally weighted cases took the most steps to converge. The two fastest algorithms were the externally weighted HALS and RHALS algorithms, taking 1.01 and 0.50 seconds, respectively.

PMF factor profiles are often visually analyzed and compared to known candidate profiles for identification. Thus, the question arises of whether or not an algorithm utilizing external weighting still contains interpretable factors in light of the decreased accuracy. Figure 7 shows the cumulative time series from one solution of the externally and internally weighted HALS algorithm, with the factors produced by the externally weighted algorithm and the internally weighted algorithm overlaid. Figure 8 shows the cumulative mass spectra of this solution, with the factors from the externally weighted algorithm and the internally weighted algorithm laid out side by side. Similar graphs for the MU and ALS algorithms are detailed in the Appendix.

Upon visual inspection, the time series produced by the externally weighted factors have peaks and troughs at almost identical times, but the magnitudes of these peaks and troughs can vary throughout the factors. Interestingly, there exists a consistent difference in the magnitude of the time series between internal and external weighting. Specifically, the external weighting algorithms seem to consistently overpredict the concentrations of the time series of the second factor compared to the internally weighted algorithms. The difference is small compared to the concentrations of the factors, but this may lead to an over interpretation of the importance of the second factor with an externally weighted algorithm. It should be noted that there exists a (similarly) large variation in relative factor signals between trials of internally weighted algorithms, although cautious researchers should be aware of the possibility that external weighting introduces extra error in this analysis. More research is needed as to why the magnitude of some factors may be over or underestimated by external weighting.

Comparing the mass spectra yields a similar analysis. Most factors of externally weighted algorithms share the same spikes of key ions as in the internally weighted ions, although sometimes at different magnitudes. Occasionally, some externally weighted mass spectra factors will look quite unsimilar to the corresponding internally weighted factor, or contain noticeable

spikes and divots in key ions, as seen in the second factor for the MU and HALS algorithm. It should also be noted that the solutions of the externally weighted algorithms could be rotated away from each other. These differences may be extreme enough to encourage users to utilize multiple trials to search for multiple solutions and eliminate bad solutions.

### 4.3.2 Comparison Between Expectation Maximization and External Weighting

515 To test the Expectation Maximization (EM) approach to uncertainties weighting as mentioned in Section 1.5.1, the weights  $\sigma_{ij}$  in the uncertainties matrix  $\Sigma$  are scaled so that  $\max_{i,j} 1/\sigma_{ij} = 1$ . Since the bulk computational component of the algorithm is constructing the matrix  $\mathbf{A}_1$  in Eq. 9, the authors in Yahaya et al. (2019) and Yahaya et al. (2021) recommend updating  $\mathbf{A}_1$  only after convergence or a maximum number of iterations of 20 or 50. They also note that applying the expectation step too early in the algorithm led to poorer performance due to the amount of error in the estimates of  $\mathbf{W}$  and  $\mathbf{H}$ . Since we wanted to  
520 apply the same stopping criterion as mentioned earlier, we chose to reconstruct  $\mathbf{A}_1$  at fixed iterations - after 10 and 20 PMF steps for different experiments. The first construction of  $\mathbf{A}_1$  is done at the 1<sup>st</sup>, 5<sup>th</sup>, and 10<sup>th</sup> step. We used the NDSVD initialization in Boutsidis and Gallopoulos (2008), and also varied the tolerance of the stopping condition between  $10^{-4}$  and  $10^{-6}$ . We summarize these results by presenting the range of average values across the different variations.

We present a comparison of external weighting and the Expectation Maximization algorithm in Table 3, using the ranges  
525 of values from the different experiments listed above. Each value is an average over 20 trials. We compare the convergence time of the algorithm, the number of steps, the weighted residual error, and the similarity to the PMF2 solution for both  $\mathbf{W}$  (correlation coefficient) and  $\mathbf{H}$  (cosine similarity).

As demonstrated in Table 3, some variations of the EM algorithm were able to outperform externally weighted HALS and RHALS in total time, as well as in weighted error. For instance, running the expectation maximization algorithm with the first  
530 calculation of  $\mathbf{A}_1$  taking place at the 10<sup>th</sup> step, and recalculating  $\mathbf{A}_1$  after 20 additional steps, until the convergence criterion with a tolerance of  $10^{-5}$  was reached, yielded an average weighted error of  $6.92 \times 10^3$  in 0.6703 seconds. This computational speed compares to RHALS, while the accuracy bests both RHALS and externally weighted HALS. However, no expectation maximization algorithm was as successful at recreating the PMF2 time series factors, as seen in column 5 of Table 3. Externally weighted RHALS and HALS also provided mass spectra factors with a higher similarity to the PMF2 factors, with the exception  
535 of two runs of the expectation maximization algorithm, one recalculating  $\mathbf{A}_1$  after 20 steps, with the first calculation at the 5<sup>th</sup> step, and the other after 10 steps, with the first calculation at the 1<sup>st</sup> step. These yielded average similarities of 0.8856 and 0.8861, respectively, although they both took over 2.90 seconds to run, much slower than any other algorithm tested besides internally weighted HALS.

We also tested how well each algorithm produced factors that were within a rotation of the PMF2 factors. As detailed in  
540 Section 1.4, the factor profiles of  $\hat{\mathbf{W}}\hat{\mathbf{H}} = (\mathbf{W}\mathbf{T}^{-1})(\mathbf{T}\mathbf{H})$  for a square matrix  $\mathbf{T}$  may be closer to the desired solution than the original factors  $\mathbf{W}$  and  $\mathbf{H}$ . To see the extent that  $\mathbf{W}$  can be rotated towards  $\mathbf{W}_{\text{PMF2}}$ , we find the matrix  $\mathbf{T}$  that minimizes the total squared difference between  $\mathbf{H}_{\text{PMF2}}$  and  $\mathbf{T}\mathbf{H}$ , and then find the average correlation between the rows of  $\mathbf{W}_{\text{PMF2}}$  and  $\mathbf{W}\mathbf{T}^{-1}$ . A symmetrical approach can be made to find the cosine similarity for  $\mathbf{H}$ . For internally weighted HALS, externally weighted HALS, and RHALS, we found average post rotation time series correlations to be 0.9391, 0.9021, and 0.8902, respectively.

545 and post rotation mass spectra similarities to be 0.9602, 0.9433, and 0.9518, respectively. When this approach was tested on HALS and RHALS using the EM algorithm, average post rotation time series correlations varied within 0.8528-0.8819 and 0.8314-0.8544, respectively, and post rotation mass spectra similarities within 0.9230-0.9363 and 0.9147-0.9287, respectively.

550 Ultimately, we found that external weighting recreated the PMF2 factors more consistently than expectation maximization, both before and after rotation. This may be due to the fact that the scaling of the weights used for the EM step is not perfectly analogous to creating a set of weights that represent the confidence of each data point. Thus, the EM method using this scaling may not capture the key error weighted patterns in the data as well as external weighting.

### 4.3.3 Complete Analysis of RHALS Algorithm

Figure 9 shows the weighted residual error for a six factor solution (equivalent to the number of factors in PMF2) over 20  
555 different non-negative SVD initializations. The average weighted residual error is  $7.2743 \times 10^3$ , with a convergence time of 0.5163 seconds over 38.45 steps per trial. If the solutions differ significantly, weighted residual error is a useful metric to pick a solution from 20 trials, and the 11<sup>th</sup> and 13<sup>th</sup> trials look like promising solutions, with weighted residual errors of  $7.1474 \times 10^3$  and  $7.1026 \times 10^3$ , respectively.

The average similarities between RHALS factors and PMF2 factors are presented in ~~FigSure~~ Figure 10, with the same  
560 similarity metrics as listed above.

We see that different initializations can lead to different solutions, in terms of both similarity to the given PMF2 solution and weighted error, suggesting that convergence to a global minima isn't always achieved. This further emphasizes the importance of using multiple initializations in order to find an optimal solution.

Only the 2<sup>nd</sup>, 11<sup>th</sup>, 13<sup>th</sup> and 17<sup>th</sup> trials have an average cosine similarity between mass spectra factors over .95, along with a  
565 correlation over .90 between time series factors. The 13<sup>th</sup> trial holds the highest similarity scores, with a time series correlation of 0.9468 and a cosine similarity of 0.9787. As seen in Figure 9, these solutions have a small weighted residual error, further justifying picking a solution with low weighted error.

How often does RHALS find a “good” solution? As we saw with the small test case, every RHALS solution contained factors very similar to the true factors. But only four out of twenty trials produced solutions surpassing 0.95 in mass spectra  
570 similarity and 0.90 in time series correlation for the entire dataset. Running 100 trials, it was found that RHALS found a solution matching these criteria 27% of the time. Thus the probability of not finding a good solution in 10 trials would be around  $(1 - .27)^{10} \approx 4.3\%$ , and  $\approx 0.2\%$  in 20 trials. This rate could vary depending on the number of factors and the ~~problem~~ itself, and should be investigated further rotational ambiguity of the problem, as in the “bad” trials RHALS may simply be finding a rotated version of the “true” solution.

575 We now explore another way one could pick a solution out of 20 trials. Researchers might be interested in seeing which factor profiles show up the most. Thus a solution that contains factors that are repeatedly found might be preferable. Thus in Figure 11 we compare a given trial to all other trials by plotting the average similarity between that trial and the other trial.

Again, cosine similarity is used for the mass spectra and the correlation coefficient is used for the time series. The bottom plot shows the two graphs averaged to give a total metric of the similarity of a solution to the other solutions.

580 Analyzing these figures, it appears that most solutions have about the same similarity to each other, with a time series correlation of around 0.87 and a mass spectra similarity of around 0.93. A few solutions, such as those in the trials between 6 and 9, can be ruled out as outliers, due to a low similarity to the other solutions. While the 12<sup>th</sup> trial yields the highest similarity between it and the other solutions (0.9148 averaged between the time series correlation and the mass spectra similarity), no solution drastically outperforms many of the other solutions. Note that the solutions from the 2<sup>nd</sup>, 11<sup>th</sup>, 13<sup>th</sup> and 17<sup>th</sup> trials  
585 all perform well in this analysis, with average similarities of 0.9028, 0.8993, 0.9016, and 0.9036, respectively.

From visual inspection of figures 9 and 11, the solution from the 13<sup>th</sup> trial is barely outperformed in the above analysis, while holding a clear lead in accuracy. Thus, it seems to be a rational choice to decide that this solution is the “best”. The solution from the 13<sup>th</sup> trial is plotted in Figures 12 and 13. Figure 12 shows the cumulative mass spectra profiles from the RHALS solution, while Figure 13 shows the cumulative time series profiles from the RHALS solution compared with those  
590 from the PMF2 solution. Again, the profiles are plotted as a sum over the entire time series and mass spectra, respectively.

The similarities of the RHALS factors to the “true” PMF2 factors are very promising; however, the magnitude of some factors differ. Specifically, RHALS underestimates the significance of the first factor, while overestimating the significance of the second and fifth factors. Again, this additional error should be considered when interpreting RHALS factors.

Another potential concern from the RHALS algorithm is the factors having a general bias towards higher or lower magni-  
595 tudes. Specifically, since RHALS enforces non negativity merely by setting negative elements to zero, in addition to regularization, it could be hypothesized that RHALS would produce factors of lower magnitude than the actual data. In Figure A1, the time series of the six RHALS factors, stacked on top of each other, are plotted against the total time series of the data matrix. Some times [we](#) see larger magnitudes from the RHALS algorithm, along with other times seeing smaller magnitudes, but there seems to be no general pattern of bias towards solutions of greater or lesser magnitude generated by RHALS.

600 Finally, we determine if a solution with a different number of factors is more optimal with the RHALS algorithm for this dataset. Figure A6 shows the convergence of weighted residual error of RHALS as the number of factors increases, with the random initialization from the 13<sup>th</sup> trial. Error dramatically decreases as the factors increase from one to six, while barely improving for solutions with a larger number of factors. Thus, a six factor solution is justified by the RHALS algorithm.

#### 4.4 Testing Rotations

605 For rotations of entire factors, we first tested the rotational algorithm detailed in [\(Paatero and Hopke \(2009\)\)](#) [Paatero and Hopke \(2009\)](#) to see whether it was applicable to RHALS. Unfortunately the pulling equations used in the Multilinear Engine (ME) do not seem to also work for RHALS.

Implementing the approach laid out in Section 2.2 to some of the solutions, we find that this method can potentially find better factorizations in the solution space. Ranging the values of  $a$  and  $b$  in Eq. (28) and Eq. (29) between 0 and 500 (which is  
610 of similar magnitude of the optimal amount of L2 regularization), we test the rotational method on the RHALS solutions from the 11<sup>th</sup>, 13<sup>th</sup> and 17<sup>th</sup> trials. The graphs are plotted in which a positive pulling value corresponds to a “pull up” of  $\mathbf{W}$  and a

“pull down” of  $\mathbf{H}$ , and vice versa for negative pulling values. The values of  $a$  and  $b$  are the magnitude of the pulling parameter. We present the rotation of the 13<sup>th</sup> solution in Figures 14 and 15, with the rotations of the solutions from other trials in the Appendix.

615 As one can see by analyzing the rotation of the solution from the three different trials, solutions that have lower weighted residual error and are closer to the target solution can be found through this simple pulling method, regardless of the direction of the pull. Interestingly, weighted residual error was able to be decreased to under 7000 when rotating the solution from the 13<sup>th</sup> trial towards larger values of  $W$  with a pulling parameter of 400. There also isn’t always an increase in values for the factors that are being pulled up, as seen in Figure A11 in the Appendix. Sometimes, in fact, as the factors are being pulled one  
620 direction, as most clearly seen with the rotated factors from the 11<sup>th</sup> trial, the matrix norms respond in the opposite direction.

The magnitude of change in the similarity to the PMF2 factors is small for all solutions, and rotations may only be worthwhile to look at when no good solution exists from the onset. The benefit of the rotations is more apparent when rotating poor solutions.

## 5 Conclusions

625 As the size of datasets has grown, computational costs have become increasingly expensive for traditional PMF algorithms. Thus, randomized and hierarchical algorithms are attractive alternative methods. Specifically, the RHALS algorithm was shown to provide a reduction in run time compared to the multiplicative update algorithm, as well as the deterministic HALS algorithm. Furthermore, we proposed a novel approach to handling uncertainties in a weighted factorization problem. While this approach, coupled with randomization, slightly reduced the accuracy of the algorithm, it dramatically decreased the computational cost. Ultimately, we showed that our weighted RHALS algorithm was able to almost completely recreate the factors in  
630 both a formed test matrix and a real dataset, and is a useful tool to finding non-negative factors of large datasets.

*Code and data availability.*

The source code for the algorithms detailed in this paper are given at <https://zenodo.org/badge/latestdoi/537111580>. The datasets used for this analysis are given at <https://doi.org/10.25810/0829-t279>.

## 635 Appendix A: ALS Derivation

The cost function  $Q$  is defined as

$$Q = \left\| \frac{\mathbf{A} - \mathbf{W}\mathbf{H}}{\mathbf{\Sigma}} \right\|_F^2 + \sum_{i=1}^m \alpha \|\mathbf{W}(i, :)\|_1 + \sum_{j=1}^n \beta \|\mathbf{H}(:, j)\|_1 + \gamma \|\mathbf{W}\|_F^2 + \delta \|\mathbf{H}\|_F^2 \quad (\text{A1})$$

where  $\alpha, \beta, \gamma$ , and  $\delta$  are regularization parameters. As with the derivation of HALS, the elementwise division of  $\Sigma$  is eliminated by considering a row or column of the residual at a time. Thus, we find

$$0 \leftarrow \frac{\partial Q^i}{\partial \mathbf{W}(i,:)} = \frac{\partial}{\partial \mathbf{W}(i,:)} (\|\mathbf{A}(i,:) - \mathbf{W}(i,:) \mathbf{H} \Sigma_i^{-1}\|_F^2 + \alpha \|\mathbf{W}(i,:)\|_1 + \sum_{j=1}^n \beta \|\mathbf{H}(:,j)\|_1 + \gamma \|\mathbf{W}(i,:)\|_2^2 + \delta \|\mathbf{H}\|_F^2) \quad (\text{A2})$$

$$0 \leftarrow \frac{\partial Q^j}{\partial \mathbf{H}(:,j)} = \frac{\partial}{\partial \mathbf{H}(:,j)} (\|\Sigma_j^{-1}(\mathbf{A}(:,j) - \mathbf{W} \mathbf{H}(:,j))\|_F^2 + \sum_{i=1}^m \alpha \|\mathbf{W}(i,:)\|_1 + \beta \|\mathbf{H}(:,j)\|_1 + \gamma \|\mathbf{W}\|_F^2 + \delta \|\mathbf{H}(:,j)\|_2^2) \quad (\text{A3})$$

where  $\Sigma_i$  is a diagonal  $n \times n$  matrix with the diagonal values equal to the  $i$ th row of  $\Sigma$ , and  $\Sigma_j$  is a diagonal  $m \times m$  matrix with the diagonal values equal to the  $j$ th column of  $\Sigma$ . Using the fact that  $\|\mathbf{X}\|_F^2 = \text{Tr}(\mathbf{X}^T \mathbf{X})$ , with  $\text{Tr}$  being the trace of the matrix, Eq. (A2) and Eq. (A3) can be rewritten as

$$0 = \frac{\partial}{\partial \mathbf{W}_i} (\text{Tr}(\Sigma_i^{-1} \mathbf{A}_i^T \mathbf{A}_i \Sigma_i^{-1} - 2 \Sigma_i^{-1} \mathbf{A}_i^T \mathbf{W}_i \mathbf{H} \Sigma_i^{-1} + \Sigma_i^{-1} \mathbf{H}^T \mathbf{W}_i^T \mathbf{W}_i \mathbf{H} \Sigma_i^{-1}) + \alpha \|\mathbf{W}(i,:)\|_1 + \sum_{j=1}^n \beta \|\mathbf{H}(:,j)\|_1 + \gamma \|\mathbf{W}(i,:)\|_2^2 + \delta \|\mathbf{H}\|_F^2) \quad (\text{A4})$$

$$0 = \frac{\partial}{\partial \mathbf{H}_j} (\text{Tr}(\mathbf{A}_j^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{A}_j - 2 \mathbf{H}_j^T \mathbf{W}^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{A}_j + \mathbf{H}_j^T \mathbf{W}^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{W} \mathbf{H}_j) + \sum_{i=1}^m \alpha \|\mathbf{W}(i,:)\|_1 + \beta \|\mathbf{H}(:,j)\|_1 + \gamma \|\mathbf{W}\|_F^2 + \delta \|\mathbf{H}(:,j)\|_2^2) \quad (\text{A5})$$

where  $\mathbf{A}(i,:) = \mathbf{A}_i$ ,  $\mathbf{A}(:,j) = \mathbf{A}_j$ ,  $\mathbf{W}(i,:) = \mathbf{W}_i$ , and  $\mathbf{H}(:,j) = \mathbf{H}_j$ . Using vector derivative rules and the fact that  $\text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{CAB}) = \text{Tr}(\mathbf{BCA})$  when  $\mathbf{ABC}$  is square, the following gradients are found.

$$0 \leftarrow -2 \mathbf{A}_i \Sigma_i^{-1} \Sigma_i^{-1} \mathbf{H}^T + 2 \mathbf{W}_i \mathbf{H} \Sigma_i^{-1} \Sigma_i^{-1} \mathbf{H}^T + \alpha \vec{1}_{1 \times k} + 2 \gamma \mathbf{W}_i \quad (\text{A6})$$

$$0 \leftarrow -2 \mathbf{W}^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{A}_j + 2 \mathbf{W}^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{W} \mathbf{H}_j + \beta \vec{1}_{k \times 1} + 2 \delta \mathbf{H}_j \quad (\text{A7})$$

In Eq. (A6) and Eq. (A7),  $\vec{1}$  is a vector of ones. Finally, the following update rules are found.

$$\mathbf{W}_i \leftarrow [(\mathbf{A}_i \Sigma_i^{-1} \Sigma_i^{-1} \mathbf{H}^T - \frac{\alpha}{2} \vec{1}_{1 \times k})(\mathbf{H} \Sigma_i^{-1} \Sigma_i^{-1} \mathbf{H}^T + \gamma \mathbf{I}_{k \times k})^{-1}]_+ \quad (\text{A8})$$

$$\mathbf{H}_j \leftarrow [(\mathbf{W}^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{W} + \delta \mathbf{I}_{k \times k})^{-1}(\mathbf{W}^T \Sigma_j^{-1} \Sigma_j^{-1} \mathbf{A}_j - \frac{\beta}{2} \vec{1}_{k \times 1})]_+ \quad (\text{A9})$$

where  $\mathbf{I}$  is the identity matrix and  $[v]_+$  returns all negative entries in  $v$  as 0.

## Appendix B: Mathematical Equivalence of Pseudoinverse Update and Ordinary Least Squares

For simplicity, we label  $(\tilde{\mathbf{W}}\tilde{\mathbf{H}}) \odot \Sigma$  as  $\tilde{\mathbf{y}}$ . Let's consider the ordinary least squares update for ~~H~~W:

$$665 \quad \mathbf{W} = \tilde{\mathbf{y}}\mathbf{H}^T(\mathbf{H}\mathbf{H}^T)^{-1}$$

$\mathbf{H}$  is a  $k \times n$  matrix that is assumed to be of full column rank, and thus can be decomposed into the exact SVD factorization  $\mathbf{H} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ , where  $\mathbf{U}$  is a  $k \times k$  orthogonal matrix,  $\mathbf{S}$  is a  $k \times k$  diagonal matrix, and  $\mathbf{V}$  is a  $n \times k$  matrix with orthonormal columns. The pseudoinverse of  $\mathbf{H}$  is defined as  $\mathbf{V}\mathbf{S}^{-1}\mathbf{U}^T$ .  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$ , and  $\mathbf{U}^T = \mathbf{U}^{-1}$ . The update for  $\mathbf{W}$  can be rewritten as

$$\begin{aligned} \mathbf{W} &= \tilde{\mathbf{y}}\mathbf{V}\mathbf{S}\mathbf{U}^T(\mathbf{U}\mathbf{S}\mathbf{U}^T)^{-1} \\ &= \tilde{\mathbf{y}}\mathbf{V}\mathbf{S}\mathbf{U}^T(\mathbf{U}\mathbf{S}^2\mathbf{U}^T)^{-1} \\ &= \tilde{\mathbf{y}}\mathbf{V}\mathbf{S}\mathbf{U}^T(\mathbf{U}\mathbf{S}^{-2}\mathbf{U}^T) \\ &= \tilde{\mathbf{y}}\mathbf{V}\mathbf{S}\mathbf{S}^{-2}\mathbf{U}^T \\ &= \tilde{\mathbf{y}}(\mathbf{V}\mathbf{S}^{-1}\mathbf{U}^T) \\ &= \tilde{\mathbf{y}}\mathbf{H}^\dagger \end{aligned} \tag{A1}$$

670 which is equal to the pseudoinverse update described in Section 2.3. A similar argument can be made for the update for  $\mathbf{H}$ .

*Author contributions.*

B.C.S. wrote the manuscript, conducted all of the numerical experiments, and derived the equations. D.K.H. advised the numerical and theoretical approach. D.K.H., M. C., ~~and H.S.~~, and J.J. reviewed and edited the manuscript. M.C. and H.S. provided the AMS data set and reference PMF2 solution.

675 *Competing interests.*

The authors declare that they have no conflict of interest.

*Acknowledgements.* B.C.S. acknowledges funding from CU SPUR (Summer Program for Undergraduate Research), D.K.H. acknowledges support from NASA 80NSSC20K0214.



## References

- Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., and Bengio, Y.: Theano: A CPU and GPU Math Compiler in Python, pp. 18–24, Austin, Texas, <https://doi.org/10.25080/Majora-92bf1922-003>, 2010.
- Boutsidis, C. and Gallopoulos, E.: SVD based initialization: A head start for nonnegative matrix factorization, *Pattern Recognition*, 41, 1350–1362, <https://doi.org/10.1016/j.patcog.2007.09.010>, 2008.
- Burred, J.: Detailed derivation of multiplicative update rules for NMF, <https://www.semanticscholar.org/paper/Detailed-derivation-of-multiplicative-update-rules-Burred/3376b4df752f2428c451e530f9c6e0ce3a3f05e4>, 2017.
- Burred, J. J.: Detailed derivation of multiplicative update rules for NMF, p. 8.
- Cichocki, A. and Phan, A.-H.: Fast Local Algorithms for Large Scale Nonnegative Matrix and Tensor Factorizations, *IEICE TRANSACTIONS on Fundamentals of Electronics, Communications and Computer Sciences*, E92-A, 708–721, [https://search.ieice.org/bin/summary.php?id=e92-a\\_3\\_708&category=A&year=2009&lang=E&abst=](https://search.ieice.org/bin/summary.php?id=e92-a_3_708&category=A&year=2009&lang=E&abst=), publisher: The Institute of Electronics, Information and Communication Engineers, 2009.
- Eckart, C. and Young, G.: The approximation of one matrix by another of lower rank, *Psychometrika*, 1, 211–218, <https://doi.org/10.1007/BF02288367>, 1936.
- Erichson, N. B., Mendible, A., Wihlbom, S., and Kutz, J. N.: Randomized Nonnegative Matrix Factorization, *Pattern Recognition Letters*, 104, 1–7, <https://doi.org/10.1016/j.patrec.2018.01.007>, arXiv:1711.02037 [cs, stat], 2018.
- Feng, X., Yu, W., and Li, Y.: Faster Matrix Completion Using Randomized SVD, in: 2018 IEEE 30th International Conference on Tools with Artificial Intelligence (ICTAI), pp. 608–615, <https://doi.org/10.1109/ICTAI.2018.00098>, iISSN: 2375-0197, 2018.
- Fornasier, M., Rauhut, H., and Ward, R.: Low-rank Matrix Recovery via Iteratively Reweighted Least Squares Minimization, *SIAM Journal on Optimization*, 21, 1614–1640, <https://doi.org/10.1137/100811404>, publisher: Society for Industrial and Applied Mathematics, 2011.
- Gillis, N.: Chapter 8: Iterative algorithms for NMF, in: *Nonnegative Matrix Factorization, Data Science*, pp. 261–305, Society for Industrial and Applied Mathematics, <https://doi.org/10.1137/1.9781611976410.ch8>, 2020.
- Halko, N., Martinsson, P. G., and Tropp, J. A.: Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions, *SIAM Review*, 53, 217–288, <https://doi.org/10.1137/090771806>, publisher: Society for Industrial and Applied Mathematics, 2011.
- Hansen, P. C. and O’Leary, D. P.: The Use of the L-Curve in the Regularization of Discrete Ill-Posed Problems, *SIAM Journal on Scientific Computing*, 14, 1487–1503, <https://doi.org/10.1137/0914086>, publisher: Society for Industrial and Applied Mathematics, 1993.
- Ho, N.-D.: Non-negative Matrix Factorization Algorithms and Applications, Ph.D. thesis, UNIVERSITÉ CATHOLIQUE DE LOUVAIN, Belgium, <https://doi.org/10.1002/9781118679852.ch15>, 2014.
- Hu, Y., Zhang, D., Ye, J., Li, X., and He, X.: Fast and Accurate Matrix Completion via Truncated Nuclear Norm Regularization, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35, 2117–2130, <https://doi.org/10.1109/TPAMI.2012.271>, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence, 2013.
- Jia, L., Liang, Y., Li, X., Lu, L., and Yan, S.: Enabling Efficient Fast Convolution Algorithms on GPUs via MegaKernels, *IEEE Transactions on Computers*, 69, 986–997, <https://doi.org/10.1109/TC.2020.2973144>, conference Name: IEEE Transactions on Computers, 2020.
- Kaloorazi, M. F. and Chen, J.: Randomized Truncated Pivoted QLP Factorization for Low-Rank Matrix Recovery, *IEEE Signal Processing Letters*, 26, 1075–1079, <https://doi.org/10.1109/LSP.2019.2920054>, conference Name: IEEE Signal Processing Letters, 2019.

- 715 Kim, E. and Hopke, P. K.: Source Identifications of Airborne Fine Particles Using Positive Matrix Factorization and U.S. Environmental Protection Agency Positive Matrix Factorization, *Journal of the Air & Waste Management Association*, 57, 811–819, <https://doi.org/10.3155/1047-3289.57.7.811>, publisher: Taylor & Francis \_eprint: <https://doi.org/10.3155/1047-3289.57.7.811>, 2007.
- Kim, J. and Park, H.: Fast Nonnegative Matrix Factorization: An Active-Set-Like Method and Comparisons, *SIAM Journal on Scientific Computing*, 33, 3261–3281, <https://doi.org/10.1137/110821172>, publisher: Society for Industrial and Applied Mathematics, 2011.
- 720 Kumar, K.: Principal component analysis: Most favourite tool in chemometrics, *Resonance*, 22, 747–759, <https://doi.org/10.1007/s12045-017-0523-9>, 2017.
- Lahabar, S. and Narayanan, P. J.: Singular value decomposition on GPU using CUDA, in: 2009 IEEE International Symposium on Parallel & Distributed Processing, pp. 1–10, <https://doi.org/10.1109/IPDPS.2009.5161058>, iSSN: 1530-2075, 2009.
- Lee, D. D. and Seung, H. S.: Learning the parts of objects by non-negative matrix factorization, *Nature*, 401, 788–791, <https://doi.org/10.1038/44565>, number: 6755 Publisher: Nature Publishing Group, 1999.
- 725 Lin, C.-J.: Projected Gradient Methods for Nonnegative Matrix Factorization, *Neural Computation*, 19, 2756–2779, <https://doi.org/10.1162/neco.2007.19.10.2756>, conference Name: Neural Computation, 2007.
- Lu, J. and Wu, L.: Technical details and programming guide for a general two-way positive matrix factorization algorithm, *Journal of Chemometrics*, 18, 519–525, <https://doi.org/10.1002/cem.894>, \_eprint: <https://analyticalsciencejournals.onlinelibrary.wiley.com/doi/pdf/10.1002/cem.894>, 2004.
- 730 Massoli, P., Stark, H., Canagaratna, M. R., Krechmer, J. E., Xu, L., Ng, N. L., Mauldin, R. L., Yan, C., Kimmel, J., Misztal, P. K., Jimenez, J. L., Jayne, J. T., and Worsnop, D. R.: Ambient Measurements of Highly Oxidized Gas-Phase Molecules during the Southern Oxidant and Aerosol Study (SOAS) 2013, *ACS Earth and Space Chemistry*, 2, 653–672, <https://doi.org/10.1021/acsearthspacechem.8b00028>, publisher: American Chemical Society, 2018.
- 735 Nasre, R., Burtcher, M., and Pingali, K.: Morph algorithms on GPUs, in: Proceedings of the 18th ACM SIGPLAN symposium on Principles and practice of parallel programming, PPOPP '13, pp. 147–156, Association for Computing Machinery, New York, NY, USA, <https://doi.org/10.1145/2442516.2442531>, 2013.
- Paatero, P.: Least squares formulation of robust non-negative factor analysis, *Chemometrics and Intelligent Laboratory Systems*, 37, 23–35, [https://doi.org/10.1016/S0169-7439\(96\)00044-5](https://doi.org/10.1016/S0169-7439(96)00044-5), 1997.
- 740 Paatero, P. and Hopke, P. K.: Rotational tools for factor analytic models, *Journal of Chemometrics*, 23, 91–100, <https://doi.org/10.1002/cem.1197>, \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/cem.1197>, 2009.
- Paatero, P. and Tapper, U.: Analysis of different modes of factor analysis as least squares fit problems, *Chemometrics and Intelligent Laboratory Systems*, 18, 183–194, [https://doi.org/10.1016/0169-7439\(93\)80055-M](https://doi.org/10.1016/0169-7439(93)80055-M), 1993.
- Paatero, P. and Tapper, U.: Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values, *Environmetrics*, 5, 111–126, <https://doi.org/10.1002/env.3170050203>, \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/env.3170050203>, 1994.
- 745 Paatero, P., Hopke, P. K., Song, X.-H., and Ramadan, Z.: Understanding and controlling rotations in factor analytic models, *Chemometrics and Intelligent Laboratory Systems*, 60, 253–264, [https://doi.org/10.1016/S0169-7439\(01\)00200-3](https://doi.org/10.1016/S0169-7439(01)00200-3), 2002.
- Paatero, P., Hopke, P. K., Begum, B. A., and Biswas, S. K.: A graphical diagnostic method for assessing the rotation in factor analytical models of atmospheric pollution, *Atmospheric Environment*, 39, 193–201, <https://doi.org/10.1016/j.atmosenv.2004.08.018>, 2005.
- 750

- Sun, D. L. and Mazumder, R.: Non-negative matrix completion for bandwidth extension: A convex optimization approach, in: 2013 IEEE International Workshop on Machine Learning for Signal Processing (MLSP), pp. 1–6, <https://doi.org/10.1109/MLSP.2013.6661924>, iISSN: 2378-928X, 2013.
- 755 Takács, G. and Tikk, D.: Alternating least squares for personalized ranking, in: Proceedings of the sixth ACM conference on Recommender systems, RecSys '12, pp. 83–90, Association for Computing Machinery, New York, NY, USA, <https://doi.org/10.1145/2365952.2365972>, 2012.
- Tan, W., C. S. F. L. L. C. W. Z. and Cao, L.: Matrix Factorization on GPUs with Memory Optimization and Approximate Computing, 26.
- Tan, W., Chang, S., Fong, L., Li, C., Wang, Z., and Cao, L.: Matrix Factorization on GPUs with Memory Optimization and Approximate Computing, in: Proceedings of the 47th International Conference on Parallel Processing, ICPP '18, pp. 1–10, Association for Computing Machinery, New York, NY, USA, <https://doi.org/10.1145/3225058.3225096>, 2018.
- 760 Ulbrich, I. M., Canagaratna, M. R., Zhang, Q., Worsnop, D. R., and Jimenez, J. L.: Interpretation of organic components from Positive Matrix Factorization of aerosol mass spectrometric data, *Atmospheric Chemistry and Physics*, 9, 2891–2918, <https://doi.org/10.5194/acp-9-2891-2009>, publisher: Copernicus GmbH, 2009.
- Wei, S., Zheng, X., Chen, D., and Chen, C.: A hybrid approach for movie recommendation via tags and ratings, *Electronic Commerce Research and Applications*, 18, 83–94, <https://doi.org/10.1016/j.elerap.2016.01.003>, 2016.
- 765 Xie, Y.-L., Hopke, P. K., and Paatero, P.: Positive matrix factorization applied to a curve resolution problem, *Journal of Chemometrics*, 12, 357–364, [https://doi.org/10.1002/\(SICI\)1099-128X\(199811/12\)12:6<357::AID-CEM523>3.0.CO;2-S](https://doi.org/10.1002/(SICI)1099-128X(199811/12)12:6<357::AID-CEM523>3.0.CO;2-S), \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/%28SICI%291099-128X%28199811/12%2912%3A6%3C357%3A%3AAID-CEM523%3E3.0.CO%3B2-S>, 1998.
- 770 Yahaya, F.: Compressive informed (semi-)non-negative matrix factorization methods for incomplete and large-scale data : with application to mobile crowd-sensing data, Ph.D. thesis, 2021.
- Yahaya, F., Puigt, M., Delmaire, G., and Roussel, G.: How to Apply Random Projections to Nonnegative Matrix Factorization with Missing Entries?, in: 2019 27th European Signal Processing Conference (EUSIPCO), pp. 1–5, <https://doi.org/10.23919/EUSIPCO.2019.8903036>, iISSN: 2076-1465, 2019.
- 775 Yahaya, F., Puigt, M., Delmaire, G., and Roussel, G.: Random Projection Streams for (Weighted) Nonnegative Matrix Factorization, in: ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 3280–3284, <https://doi.org/10.1109/ICASSP39728.2021.9413496>, 2021.
- Zhang, S., Wang, W., Ford, J., and Makedon, F.: Learning from Incomplete Ratings Using Non-negative Matrix Factorization, pp. 549–553, <https://doi.org/10.1137/1.9781611972764.58>, 2006.

**Table 1.** ~~Average statistics Comparison of different algorithms over 5 different SVD initializations, tolerance= $10^{-4}$  (the number of operations per step for MU, tolerance is  $10^{-5}$  each algorithm used. We also track the number of elementwise multiplications and a random initialization divisions (implemented in Matlab as “.” and “./”) , L1 Regularization=1, L2 Regularization=50 in the third column.~~

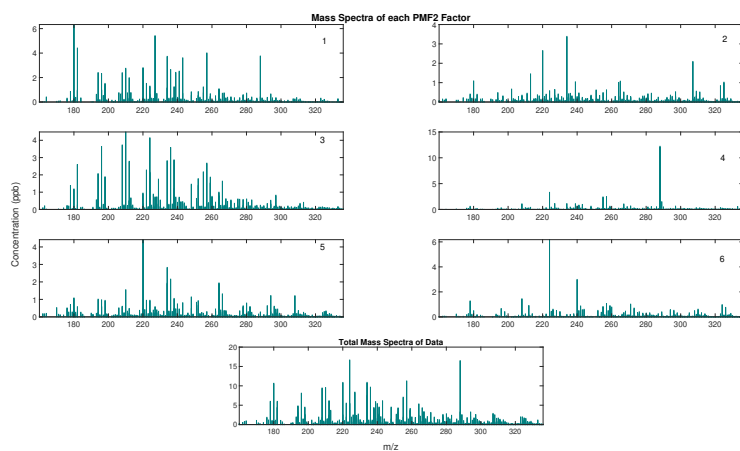
Comparison of Cost/Step for Different Algorithms		
Algorithm	Number of Operations Using Matrix Multiplication	Number of Elementwise Multiplications/Divisions
ALS (IW)	$\mathcal{O}(mnk^2)$	$2mnk$
ALS (EW)	$\mathcal{O}(mnk)$	0
MU (IW)	$\mathcal{O}(mnk)$	$4mn + 4(m+n)k$
MU (EW)	$\mathcal{O}(mnk)$	$4(m+n)k$
MU (Randomized)	$\mathcal{O}((k+l)nk)$	$4((k+l)+n)k$
HALS (IW)	$\mathcal{O}(mnk^2)$	$2mnk$
HALS (EW)	$\mathcal{O}(mnk)$	0
RHALS	$\mathcal{O}((k+l)nk)$	0
Post Processing	$\mathcal{O}(mnk)$	0

**Table 2.** Average statistics of different algorithms over 5 different SVD initializations, tolerance= $10^{-4}$  (for MU, tolerance is  $10^{-5}$  and a random initialization), L1 Regularization=1, L2 Regularization=50.

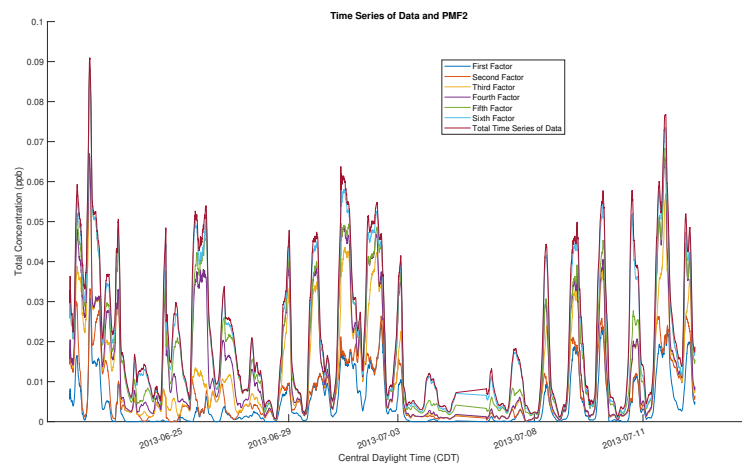
Comparison of Algorithms			
Algorithm	Total Time (s)	Steps	Weighted Error
ALS (IW)	317.08	42.2	$6.48 \times 10^3$
ALS (EW)	12.21	29.0	$7.35 \times 10^3$
MU (IW)	19.31	76.4	$6.52 \times 10^3$
MU (EW)	4.06	51.2	$7.11 \times 10^3$
MU (Randomized)	4.02	53.6	$7.76 \times 10^3$
HALS (IW)	33.87	19.4	$6.45 \times 10^3$
HALS (EW)	1.01	31.2	$7.08 \times 10^3$
RHALS	0.50	36.6	$7.15 \times 10^3$

**Table 3.** Average statistics of external weighting (EW) versus Expectation Maximization (EM) algorithms over 20 trials. Internally weighted (IW) HALS is provided as a reference. Ranges of the values of the algorithms run using different variations expectation maximization (EM) steps are presented. The correlation of the columns of **W** and similarity of the rows of **H** to the PMF2 solution are also listed.

Comparison of Algorithms					
Algorithm	Total Time (s)	Steps	Weighted Error	Correlation of <b>W</b>	Similarity of <b>H</b>
HALS (IW)	40.70	19.65	$6.46 \times 10^3$	0.8756	0.9134
HALS (EW)	1.07	31.80	$7.09 \times 10^3$	0.8414	0.8828
RHALS (EW)	0.56	38.35	$7.27 \times 10^3$	0.8475	0.9034
HALS (EM)	0.51-2.93	13.50-83.90	$6.71-7.22 \times 10^3$	0.7759-0.8306	0.8434-0.8861
RHALS (EM)	0.33-1.19	20.30-61.05	$7.11-7.54 \times 10^3$	0.7930-0.8221	0.8301-0.8694

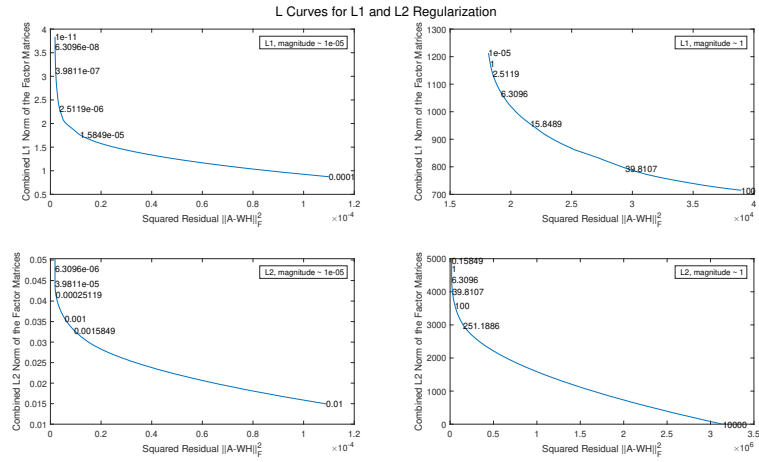


**Figure 1.** Mass spectra profiles of the 6 PMF2 factors labeled by order, as well as the total mass spectra concentration of the data.

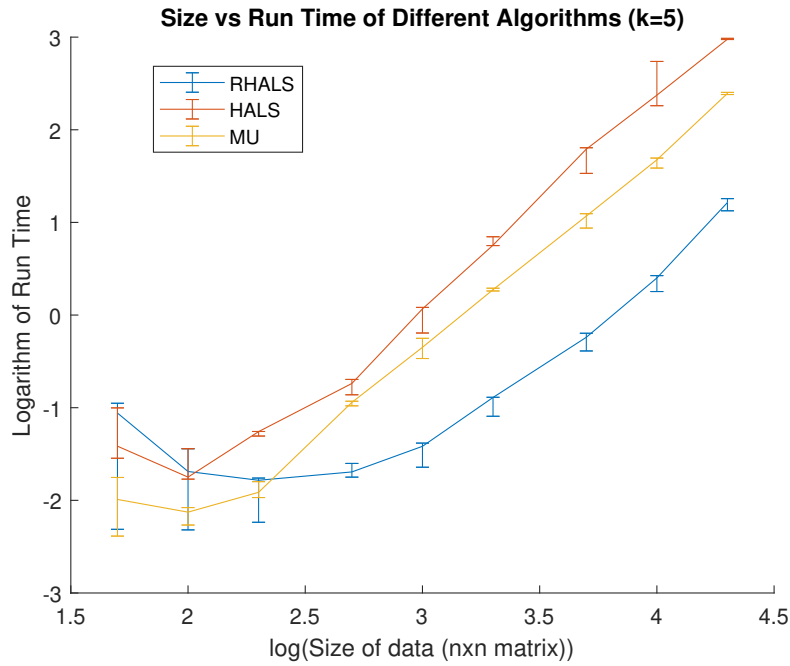


**Figure 2.** Time series profiles of the six PMF2 factors, overlaid with the total time series concentration of the data. A rolling average is used with values representing average concentrations over the previous two hours. The individual time series are stacked on top of each other in order to compare the total PMF2 time series to the overall time series of the data.

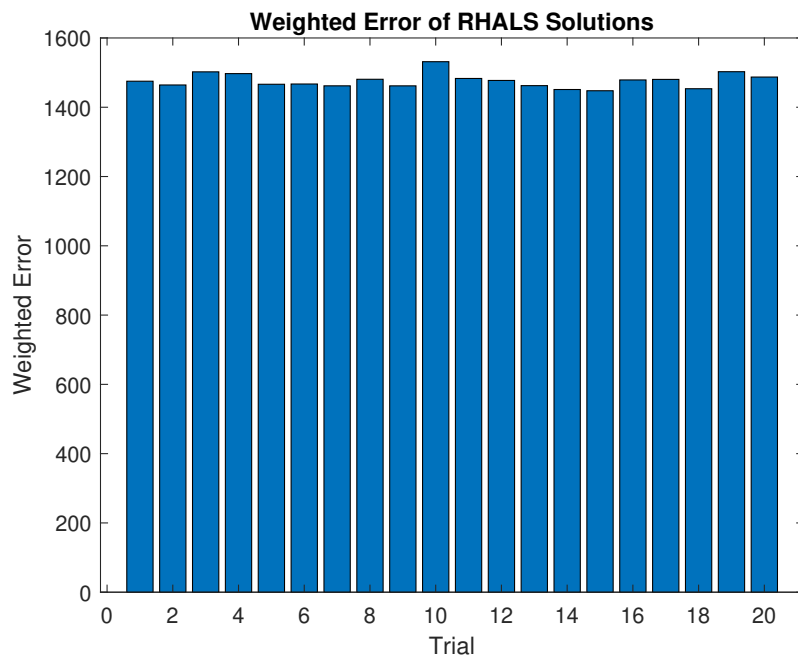




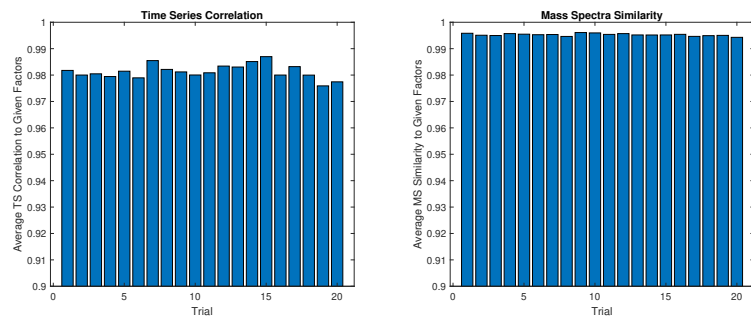
**Figure 3.** L Curves for L1 and L2 Regularization, applied to data with different average values. The top left graph shows L1 regularization for data with a magnitude around  $10^{-5}$ ; the top right graph shows L1 regularization for data with a magnitude around 1; the bottom left graph shows L2 regularization for data with a magnitude around  $10^{-5}$ ; the bottom right graph shows L2 regularization for data with a magnitude around 1.



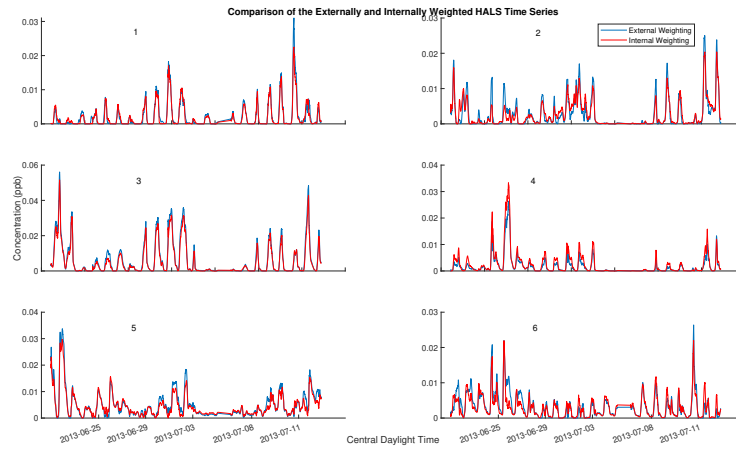
**Figure 4.** Size of data matrix vs run time (in seconds) in RHALS algorithm, performed with three different random seeds (containing the absolute value of random normal variables). The x-axis shows the base 10 log of  $n$ , and the y-axis plots the base 10 log of run time. The median run time is plotted, with the error bars plotting the maximum and minimum run times. Note that run time initially decreases due to a coincidental reduction in the number of steps to convergence, but problems of environmental interest are generally located on the right of the graph or beyond its right edge.



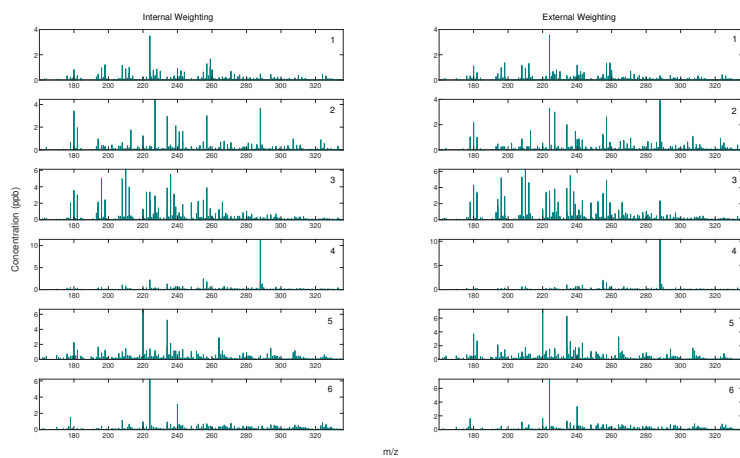
**Figure 5.** RHALS error over 20 trials. Mean: ~~1476.6~~1476.60, Standard Deviation: ~~20.4740~~20.47.



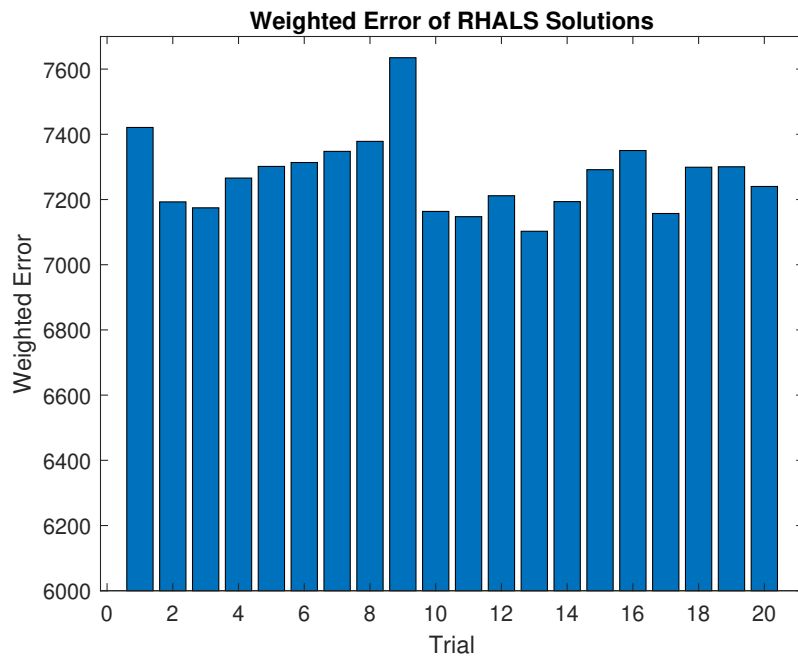
**Figure 6.** Average mass spectra and time series similarity to formed factors in the small test case.



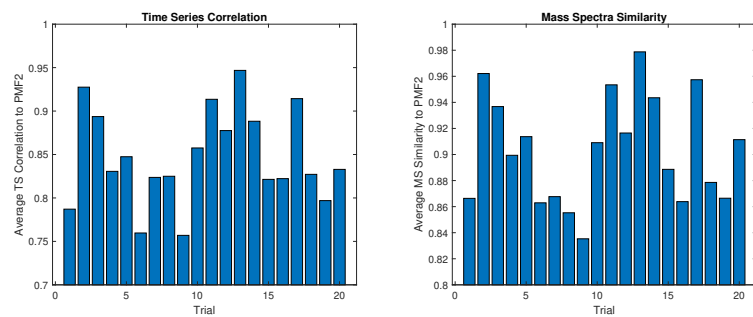
**Figure 7.** Comparison of two hour rolling average of time series for externally Weighted and internally weighted HALS factors. Externally weighted error:  $7.1120 \times 10^3$   $7.1321 \times 10^3$ . Internally weighted error:  $6.4657 \times 10^3$   $6.4841 \times 10^3$ .



**Figure 8.** Comparison of mass spectra for externally weighted and internally weighted HALS factors.

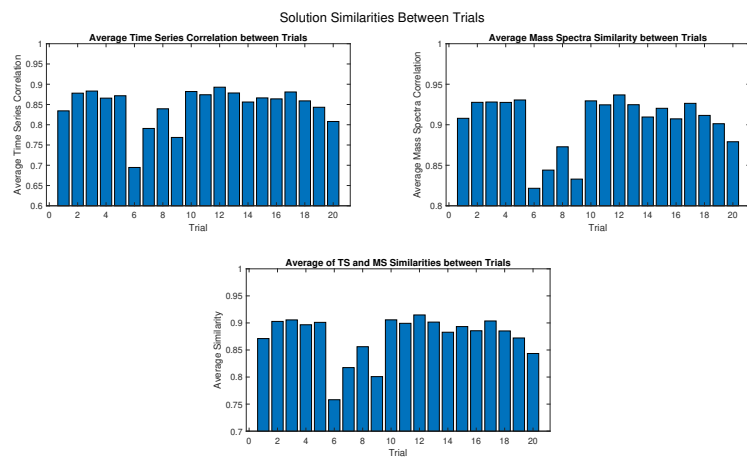


**Figure 9.** RHALS error over 20 trials. Mean Error:  $7.2743 \times 10^3$ , Standard Deviation of Error: ~~420.81~~120.8. Mean Time: 0.5163, Standard Deviation of Time: 0.0292. Mean Steps: 38.45 Standard Deviation of Steps: ~~6.8785~~6.88.

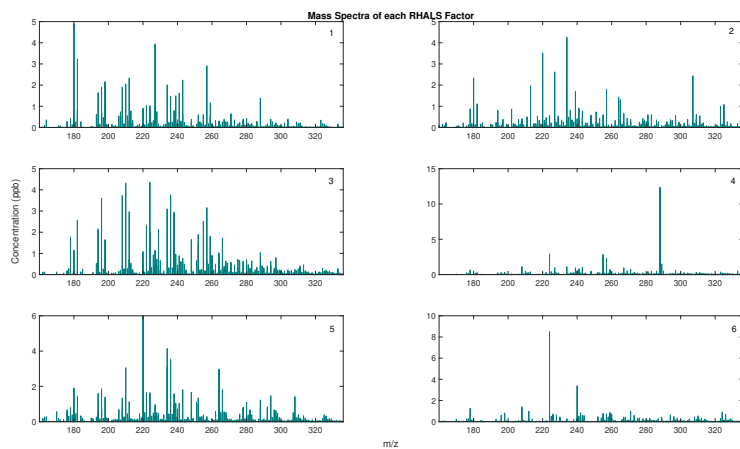


**Figure 10.** Average RHALS mass spectra and time series similarities to PMF2.

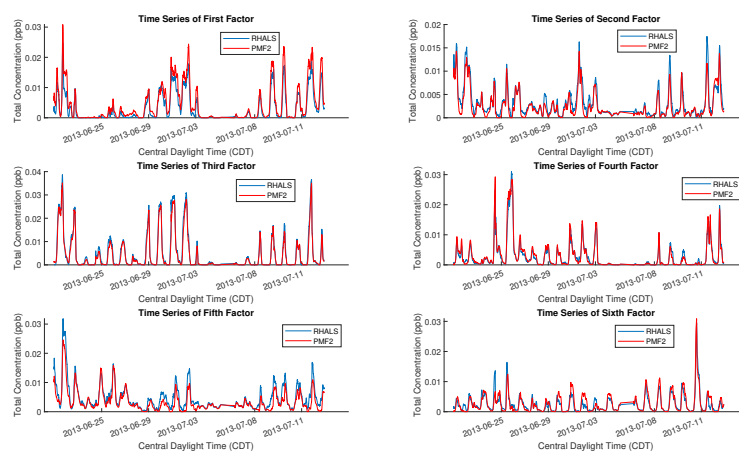




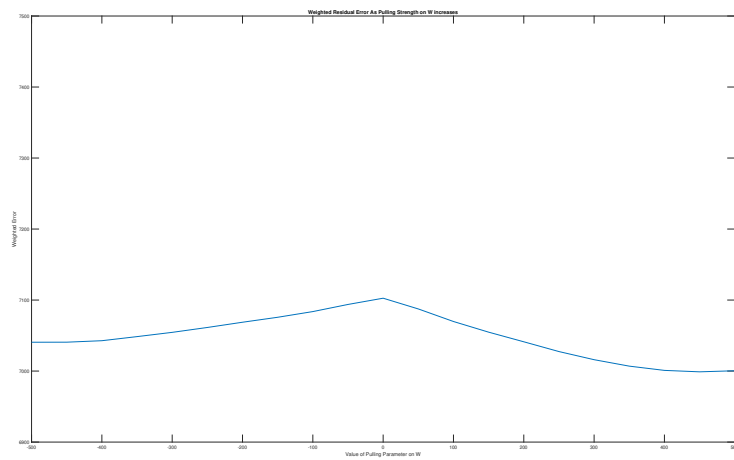
**Figure 11.** Similarity between the solutions of the different trials for RHALS



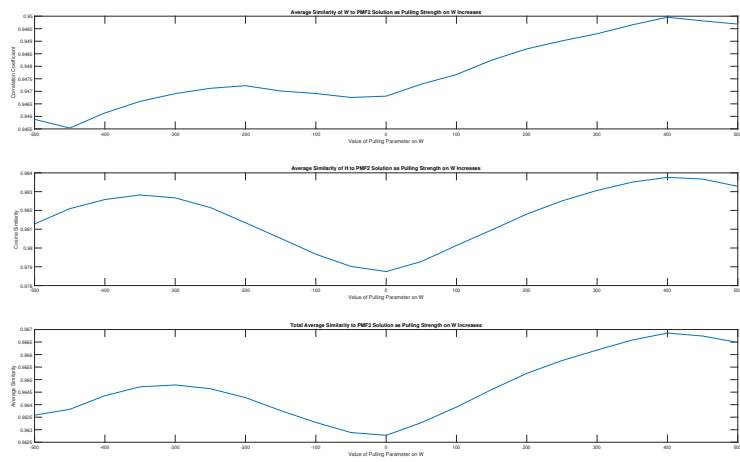
**Figure 12.** RHALS Mass Spectra Series Factors, labeled by order.; First factor similarity with PMF2=0.9631. Second factor similarity with PMF2=0.9582. Third factor similarity with PMF2=0.9912. Fourth factor similarity with PMF2=0.9970. Fifth factor similarity with PMF2=0.9795. Sixth factor similarity with PMF2=0.9835.



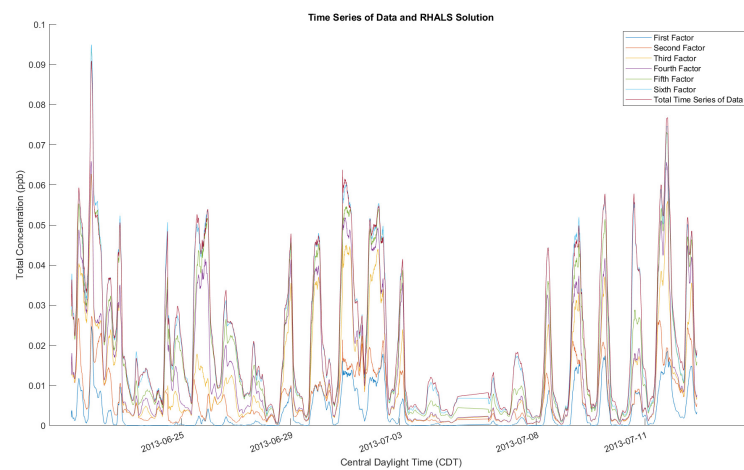
**Figure 13.** Two hour rolling average of RHALS time series factors overlaid with PMF2 factors.



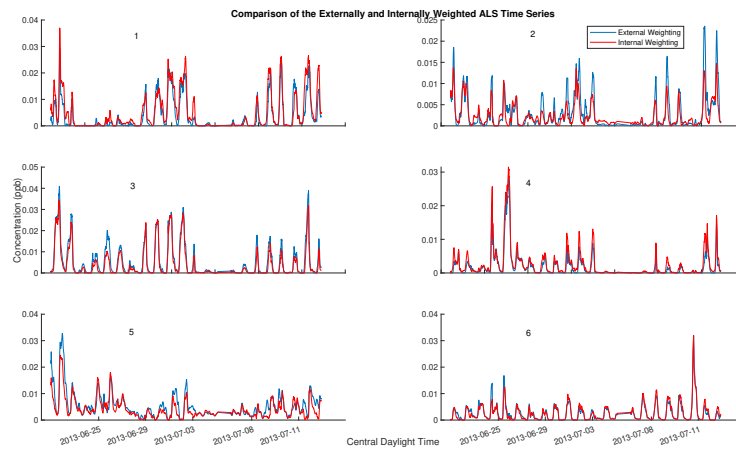
**Figure 14.** Weighted residual error of rotated solutions from the 13<sup>th</sup> trial of RHALS.



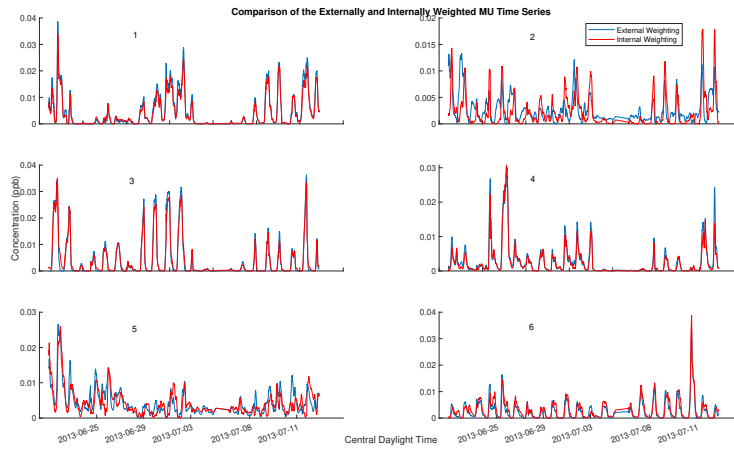
**Figure 15.** Similarities to PMF2 solution of rotated solutions from the 13<sup>th</sup> trial of RHALS.



**Figure A1.** RHALS Time Series Factors laid out with total time series from data, plotted as a two hour rolling average.

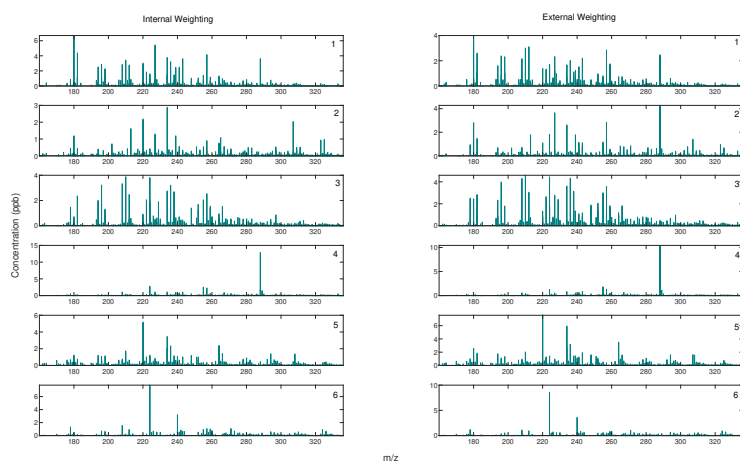


**Figure A2.** Comparison of two hour rolling average of time series for externally Weighted and internally weighted ALS factors. Externally weighted error:  $7.2465 \times 10^3$   $\underline{7.3666 \times 10^3}$ . Internally weighted error:  $6.4853 \times 10^3$   $\underline{6.4768 \times 10^3}$ .

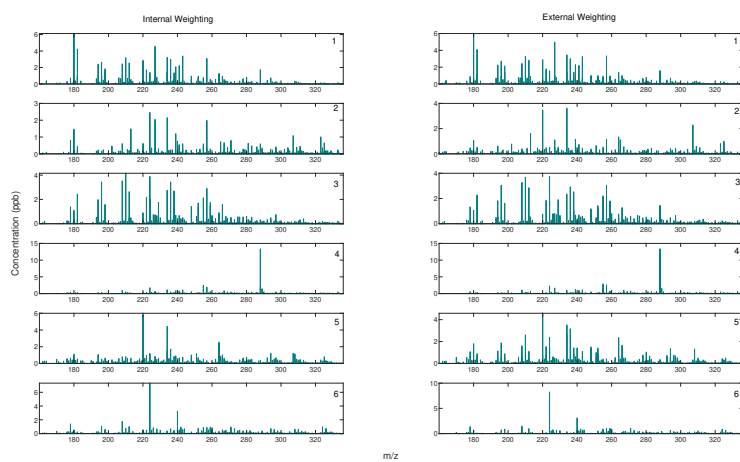


**Figure A3.** Comparison of two hour rolling average of time series for externally Weighted and internally weighted MU factors. Externally weighted error:  $7.0239 \times 10^3$   $6.9457 \times 10^3$ . Internally weighted error:  $6.3900 \times 10^3$   $6.4292 \times 10^3$ .

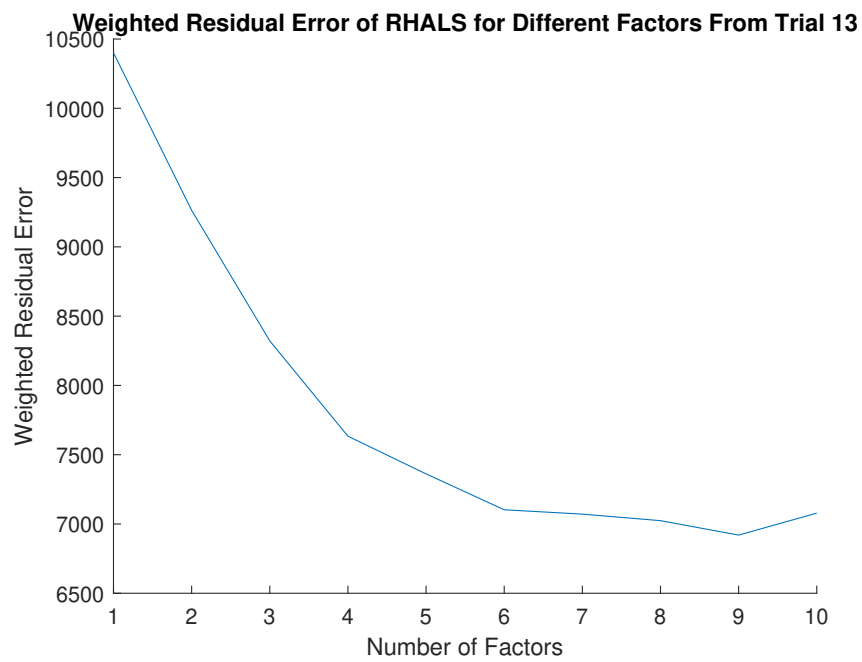




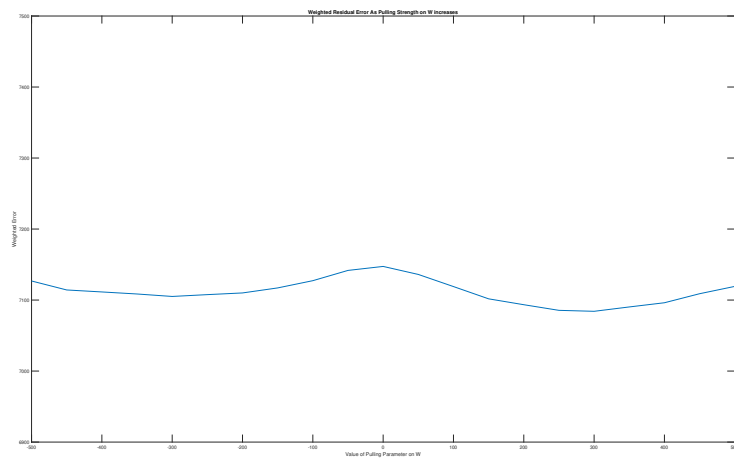
**Figure A4.** Comparison of mass spectra for externally weighted and internally weighted ALS factors.



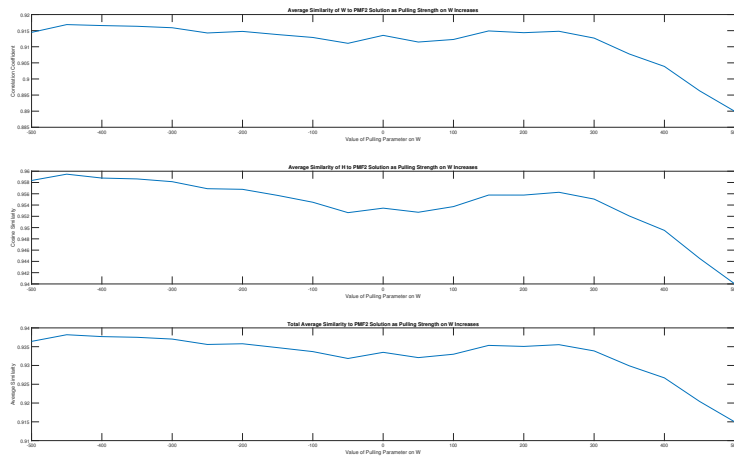
**Figure A5.** Comparison of mass spectra for externally weighted and internally weighted MU factors.



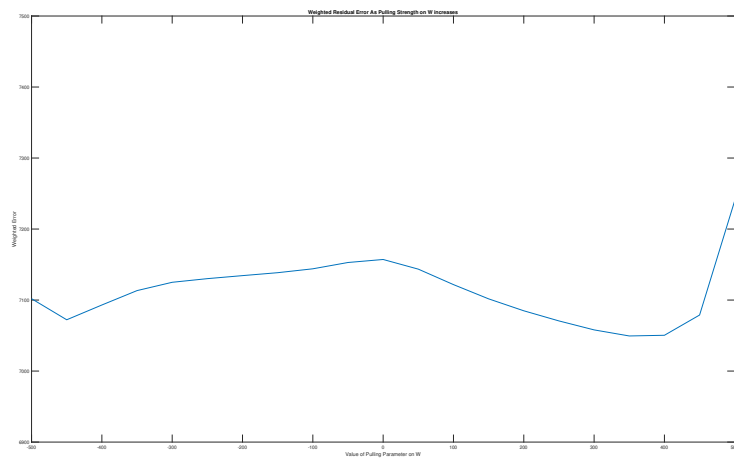
**Figure A6.** Number of factors versus weighted residual error for the RHALS algorithm.



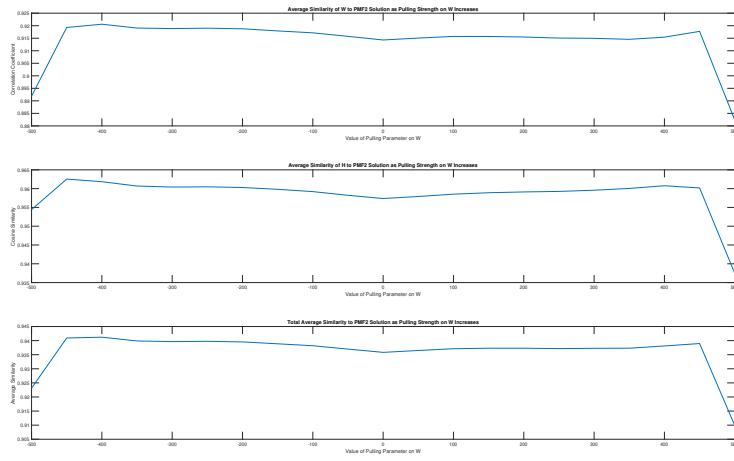
**Figure A7.** Weighted residual error of rotated solutions from the 11<sup>th</sup> trial of RHALS.



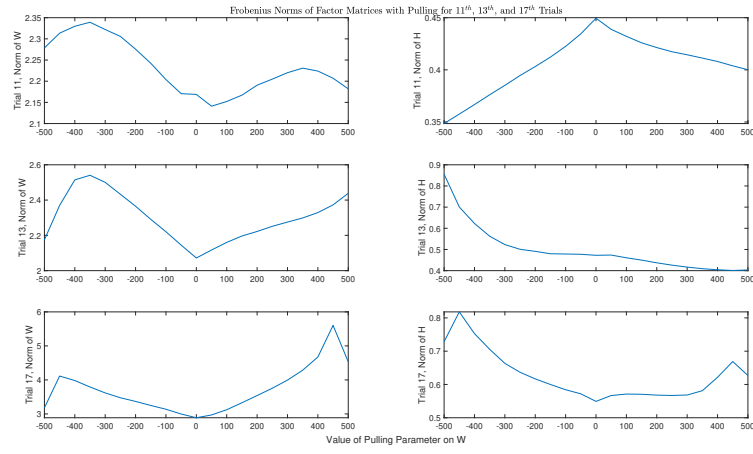
**Figure A8.** Similarities to the PMF2 solution of rotated solutions from the 11<sup>th</sup> trial of RHALS. The total similarity refers to an average of the similarity metrics for the time series and mass spectra factors.



**Figure A9.** Weighted residual error of rotated solutions from the 17<sup>th</sup> trial of RHALS.



**Figure A10.** Similarities to PMF2 solution of rotated solutions from the 17<sup>th</sup> trial of RHALS.



**Figure A11.** Frobenius norms of rotated factor matrices for the 11<sup>th</sup>, 13<sup>th</sup>, and 17<sup>th</sup> trials.