# Structured Nonnegative Matrix Factorization for Traffic Flow Estimation of Large Cloud Networks

Syed Muhammad Atif[a], Nicolas Gillis[c,*], Sameer Qazi[b], Imran Naseem[b,d]

[a]*Graduate School of Science and Engineering, Karachi Institute of Economics and Technology,Korangi Creek, Karachi, 75190, Sindh, Pakistan*

[b]*College of Engineering, Karachi Institute of Economics and Technology,Korangi Creek, Karachi, 75190, Sindh, Pakistan*

[c]*Department of Mathematics and Operational Research Faculté polytechnique, Université de Mons,Rue de Houdain 9, 7000, Mons, Belgium*

[d]*School of Electrical, Electronic and Computer Engineering, The University of Western Australia,35 Stirling Highway, Crawley, 6009, Western Australia, Australia*

## Abstract

Network traffic matrix estimation is an ill-posed linear inverse problem: it requires to estimate the unobservable origin destination traffic flows, $X$, given the observable link traffic flows, $Y$, and a binary routing matrix, $A$, which are such that $Y = AX$. This is a challenging but vital problem as accurate estimation of OD flows is required for several network management tasks. In this paper, we propose a novel model for the network traffic matrix estimation problem which maps high-dimension OD flows to low-dimension latent flows with the following three constraints: (1) nonnegativity constraint on the estimated OD flows, (2) autoregression constraint that enables the proposed model to effectively capture temporal patterns of the OD flows, and (3) orthogonality constraint that ensures the mapping between low-dimensional latent flows and the corresponding link flows to be distance preserving. The parameters of the proposed model are estimated with a training algorithm based on Nesterov accelerated gradient and generally shows fast convergence. We validate the proposed traffic flow estimation model on two real backbone IP network datasets, namely Internet2 and GÉANT. Empirical results show that the proposed model outperforms the state-of-the-art models not only in terms of tracking the individual OD flows but also in terms of standard performance metrics. The proposed model is also found to be highly scalable compared to the existing state-of-the-art approaches.

*Keywords:* network traffic matrix estimation, nonnegative matrix factorization, Nesterov accelerated gradient, autoregressive model, graph embedding, distance preserving transformation

*Corresponding author: (Nicolas Gillis) Email: nicolas.gillis@umons.ac.be., Tel.: +32-(0)65-374680, Fax: +32-(0)65-374500;

*Email addresses:* s.m.atif@kiet.edu.pk. (Syed Muhammad Atif), nicolas.gillis@umons.ac.be. (Nicolas Gillis), sameer.qazi@pafkiet.edu.pk. (Sameer Qazi), imrannaseem@pafkiet.edu.pk,imran.naseem@uwa.edu.au. (Imran Naseem)

*URL:* https://sites.google.com/site/nicolasgillis/ (Nicolas Gillis)

## 1. Introduction

Since the emergence of Internet, network parameter estimation has emerged to be an important but challenging research topic whose goal is to infer the network parameters that cannot be measured, such as end-to-end traffic volume and link delays. The significance of this research domain is due to its vital role in the network engineering tasks such as capacity planning [1], network optimization [2], expansion and monitoring of the network [3, 4], congestion avoidance [5], and anomaly detection [6]. In this paper, we focus on the network traffic matrix estimation problem that estimates the origin-to-destination traffic volumes using readily available routing matrix and easily observable link traffic volume.

Consider a network which consists of $m$ links and $n$ origin-destination (OD) pairs. Let the network be observed at time $t$, and let us denote $\boldsymbol{x}_t \in \mathbb{R}_+^n$ and $\boldsymbol{y}_t \in \mathbb{R}_+^m$ the OD flow and link flow vectors at time $t$, respectively. Further, let $A \in \{0, 1\}^{m \times n}$ be the binary routing matrix of the network such that $A(i, j) = 1$ if the routing path for the $j$th OD pair of the network passes through the $i$th link of the network. Then, the following relationship between link flow vector $\boldsymbol{y}_t$ and OD flow vector $\boldsymbol{x}_t$ holds:

$$\boldsymbol{y}_t = A\boldsymbol{x}_t. \tag{1}$$

When the network is monitored for a time period $T$, (1) becomes:

$$Y = AX \tag{2}$$

where the link flow matrix $Y \in \mathbb{R}_+^{m \times T}$ and traffic matrix $X \in \mathbb{R}_+^{n \times T}$ are formed by stacking the column vectors $\boldsymbol{y}_t$'s and $\boldsymbol{x}_t$'s, for $t = 1, 2, \ldots, T$.

The traffic estimation problem is the linear inverse problem of recovering $X$ from $A$ and $Y$ using (2). More precisely, the goal is to estimate the OD flows $\boldsymbol{x}_{T+i}$ at the future timestamps $T + i$ (for $i = 1, 2, \ldots$) provided that the corresponding link flows $\boldsymbol{y}_{T+i}$, the binary routing matrix $A$, and the historical OD flows $X \in \mathbb{R}_+^{n \times T}$ and link flows measurements $Y \in \mathbb{R}_+^{m \times T}$ are given. It is an ill-posed problem for all practical networks because the number of links in the networks is typically much smaller than the OD pairs in the network, that is, $m \ll n$. Hence, finding a solution for this problem is challenging because of the following two reasons:

(1) Due to the rapid advancement in the network technology and its abundance at an affordable price, the size of modern networks are growing day by day [7]. Hence, the difference between the number of links $m$ and that of OD pairs $n$ is becoming more and more prominent.

(2) Due to the increase in diversity of applications running on the modern Internet, the statistical features of traffic are being further complicated [8].

Dimensionality reduction techniques play vital rule in network traffic estimation problem. It is because OD flows when observed over a period of time $T$, that is, the OD flow matrix $X \in \mathbb{R}_+^{n \times T}$, have diurnal cycles hence high-dimensional OD flows can be mapped to low-dimensional latent flows using a dimensionality reduction technique such as the truncated SVD. Let us consider the following low-rank approximation

$$X \approx PQ, \tag{3}$$

where $P \in \mathbb{R}^{n \times k}$ and $Q \in \mathbb{R}^{k \times T}$ are the left and right low-rank factors, and $k < m \ll n$ is the factorization rank. The matrix $P$ encodes the spatial features of OD flows whereas $Q$ encodes the latent flows associated with the OD flow matrix $X$. Plugging (3) in (2), we obtain

$$Y \approx APQ. \tag{4}$$

For a backbone Internet network, it can be safely assumed that the routing matrix $A$ and spatial feature matrix $P$ remain stable over a long period of time. Thus, the highly ill-posed problem of estimating OD flows using (2) is transformed with the help of dimensionality reduction into the well-posed problem of estimating latent flows $Q$ using (4), given that $P$ has been estimated using historical data. In fact, given a new link flow observation $\boldsymbol{y}_t$, we can estimate $\boldsymbol{x}_t$ by solving $\min_{\boldsymbol{q}_t \in \mathbb{R}^k} \|\boldsymbol{y}_t - AP\boldsymbol{q}_t\|_2$ whose solution is given by $\boldsymbol{q}_t = (AP)^\dagger \boldsymbol{y}_t$ where $(AP)^\dagger$ is the (left) pseudoinverse of $AP \in \mathbb{R}^{m \times k}$, and obtain $\boldsymbol{x}_t = P\boldsymbol{q}_t$.

## 1.1. Related work

Several approaches have been proposed recently for traffic matrix estimation. Among them dimensionality reduction techniques are the most popular and they are the focus of this paper. The reason why dimensionality reduction techniques can be employed in network traffic estimation is because (1) the backbone Internet traffic is highly concentrated, that is, most of the OD flows use only few links of the network, and (2) OD flows have spatial and temporal similarities, that is, nearby OD pairs have similar traffic and OD flows have diurnal cycles. Soule et al. [9] proposed to employ the singular value decomposition (SVD) [10], a well-known unconstrained dimensionality reduction technique, to develop a traffic estimation model. A disadvantage of their model is that it uses the Moore-Penrose pseudoinverse hence requiring an additional step to suppress possible negative OD flows in the initial estimate of the model. In [11], Kumar et al. proposed a traffic estimation model based on a relatively new unconstrained dimensionality reduction technique known as CUR [12]. Its key benefits over SVD are mainly computational efficiency and the interpretability of the low-rank factors that are directly derived from the given data. Recently, Qazi et al. [13] proposed to use the demand matrix and the traffic probability matrix pair to transform the ill-posed traffic estimation problem into an equivalent well-posed problem. However, the two proposed methods [11, 13] suffer from the limitations similar to that of [9] due to the use of the Moore-Penrose pseudoinverse.

Some researchers have shown keen interest in different neural network architectures to design an effective model for traffic flow estimation. For example, Jiang et al. in [14] proposed to employ feed forward neural network (FFNN) for traffic flow estimation modeling. The proposed model BPTME takes link flows as input to FFNN whereas OD flows are yielded as output by the network. The network is trained by back error propagation algorithm. Zhou et al. in [15] proposed an enhancement over BPTME [14] by injecting the routing information into the FFNN as input for improved performance. Due to the architectural simplicity, FFNN does not scale well with size, complexity and dynamics of traffic flows in the modern networks, hence Nie et al [16] proposed to a traffic flow estimation model based on deep belief network (DBN). In [17], Zhao et al. used long short term memory recurrent neural network (LSTM-RNN). DLTMP [17] is found to perform better compared to the contemporary approaches in capturing spatiotemporal dependencies of traffic flows because LSTM-RNN has cyclic connections over time.

Genetic or evolutionary algorithms are also utilized by some researchers in the context of traffic flow estimation [18, 19]. These typically employ quantum-behaved particle swarm optimization (QPSO). In [20], Lu et al. used multi-fractal discrete wavelet transform (MDWT) to split traffic matrix into different frequency component then train the neural network to predict low and high frequency component of traffic matrix. Kumar et al. in [21] proposed a multi-view subspace learning technique for traffic flow estimation. They proposed a novel robust approach to obtain traffic flows from multiple traffic views yielded from rather inexpensive existing methods.

In this paper, we consider graph embedding and nonnegative matrix factorization (NMF) to perform traffic flow estimation. Graph embedding has been successfully applied in various

research fields such as computer vision and recommender systems. Interested readers may refer to the recent survey paper [22] by Cai et al. on graph embedding for further details. Emami et al. [23] have recently used graph embedding in the context of traffic flow estimation. Their approach blends graph embedding with convolution neural network (CNN) for better traffic flow estimation. Moreover, graph embedding has also been used by several researchers in combination with matrix factorization. A seminal research work in the context of NMF along with graph embedding is [24] by Cai et al. Roughan et al. [25] have used classical NMF in the scenario of network traffic flow estimation. An autoregression based approach has been proposed recently for capturing temporal dependencies in [26] by Yu et al. This autoregressive approach maps the high-dimensional time series data into a low-dimensional latent time series. To avoid overfitting, each latent timeseries is autoregressed independently. Further, the authors of [26] proved that their proposed autoregressive approach has an equivalent graph representation, thus conventional solvers used for graph embedding approaches can be used for this autoregression approach as well.

### 1.2. Contribution and outline of the paper

Motivated from the benefits of NMF and autoregression approaches, and keeping in mind the limitations of the existing traffic flow estimation models, we propose a novel model for traffic flow estimation that is developed using a unique combination of NMF, autoregression and orthogonality.

NMF is a constrained dimensionality reduction technique which ensures nonnegativity hence is the natural choice for nonnegative network traffic flow estimation. It has already been successfully applied in many engineering applications, including but not limited to, computer vision, signal processing [27], hyperspectral sensing [28], and background subtraction [29]. However to the best of our knowledge, we are the first to conduct a comprehensive research with conclusive results for the use of NMF in the context of network traffic flow estimation. Along with the nonnegativity constraint, we also impose two additional constraints on the proposed model for effectively capturing the spatial and temporal features in the traffic flows of the underlying network. The key contribution of this paper is threefold:

(1) A novel multi-constrained nonnegative matrix factorization model is proposed for network traffic flow estimation. The proposed model consists of a dimensionality reduction for capturing spatial features of the network traffic flows with a novel combination of three additional constraints: a) nonnegativity, b) autoregression, and c) orthogonality. The autoregression constraint allows a dynamic learning of autoregression weights and an effective representation of negative correlation between different OD flows which are prominent features of the proposed model and are not available in the conventional graph embedding approaches. The orthogonality constraint provides a distance preserving transformation from link flows to latent flows enabling effective clustering of the given link flows according to their temporal similarities.

(2) We propose an efficient training algorithm using the fast gradient method of Nesterov, while we use a noise resilient initialization strategy that provides a deterministic seeding point to the training algorithm [30].

(3) We illustrate the effectiveness of the proposed model and algorithm on two publicly available Internet backbone traffic data sets, namely Internet2 [31] and GÉANT [32]. The empirical

4

results show that the new model performs competitively with the current state-of-the-art models.

The rest of the paper is organized as follows. Section 2 introduces the proposed model in detail, highlighting the key components and their advantages along with differences with existing approaches. Section 3 presents our proposed algorithm to tackle our new model. In Section 4, the proposed model is extensively evaluated on two publicly available Internet traffic data sets, and we show that our proposed model competes favorably with state-of-the-art algorithms.

## 2. New model for network traffic flow estimation

Let us formulate our proposed NMF-based network traffic estimation model for a network that is monitored for time period $T$, and consists of $m$ links and $n$ OD pairs. For this time period, we are given $X \in \mathbb{R}^{n \times T}$, and our goal is to construct a model allowing us to predict $x_t$ for $t > T$, given $A$ and $y_t$; see Section 1. More precisely, we are given

- the routing matrix $A \in \mathbb{R}^{m \times n}$,

- the OD flow matrix $X \in \mathbb{R}^{n \times T}$,

- a factorization rank $k$, and

- a lag set $\mathcal{L} = \{l : l \in \mathbb{N} \text{ and } l \ll T\}$ with $L = \max(\mathcal{L})$. The lag set contains the indices $l$ indicating a dependency between the $t$th and $(t - l)$th time points, such as in autoregressive models; see Section 2.2 for more details.

The model aims at computing $(W, H, \Omega)$ by solving the following optimization problem:

$$\min_{W \in \mathbb{R}^{n \times k}, H \in \mathbb{R}^{k \times T}, \Omega \in \mathbb{R}^{k \times L}} \|X - WH\|_F^2 \tag{5a}$$

$$\text{such that } W \geq \mathbf{0}, H \geq \mathbf{0}, \Omega \geq \mathbf{0}, \tag{5b}$$

$$H(p, t) = \sum_{l \in \mathcal{L}} \Omega(p, l) H(p, t - l)$$

$$\text{for } 1 \leq p \leq k, \ L < t \leq T, \tag{5c}$$

$$\Omega(p, l) > 0,$$

$$\mathcal{A} = AW, \quad \mathcal{A}^{\top} \mathcal{A} = I_k, \tag{5d}$$

where $\|.\|_F$ denotes the Frobenius norm, and the objective (5a) is the data fitting term of the model. Let us discuss the constraints in the above model:

- (5b) ensures nonnegative entries in $W$, $H$ and $\Omega$.

- (5c) imposes that each entry $H(p, t)$ of matrix $H$ is the nonnegative weighted sum of entries in the $p$th row of matrix $H$ preceding $H(p, t)$. This is motivated by the autoregression model; see Section 2.2 for the details.

- (5d) defines a rank-$k$ compact routing matrix $\mathcal{A} = AW \in \mathbb{R}^{m \times k}$ that has orthonormal columns as $I_k$ stands for a $k \times k$ identity matrix; see Section 2.1 for the details.

5

In the context of dimensionality reduction, the low-rank matrices $W$ and $H$ are generally called the basis and embedding (encoding) matrix, respectively. However, in the context of traffic matrix estimation, $W$ contains the spatial features of the OD flow matrix $X$ in low dimension, while $H$ can be interpreted as a latent flow matrix defining the latent flows of the underlying network.

In order to incorporate the constraints (5b)-(5d), we will resort to regularization so that the constraints will not be strictly satisfied, but will tend to be. Using regularization, the optimization task in (5) can be formulated as a regularized NMF problem; see Section 3.1.

Let us now discuss the two key constraints in (5), namely orthogonality (5d) and autoregression (5c).

### 2.1. Orthogonality constraint

The constraint $\mathcal{A}^T \mathcal{A} = I_k$ can be equivalently written as $W^T(A^T A)W = I_k$, which imposes $k^2$ constraints on $W \in \mathbb{R}^{n \times k}$. Note that $\mathcal{A} = AW$ is also nonnegative since $A$ and $W$ are, and hence this orthogonality constraint requires the columns of $\mathcal{A}$ to have disjoint supports (that is, the set of nonzero entries of the columns of $AW$ do not intersect). This is related to the so-called orthogonal NMF model; see [33] and the references therein. Hence this constraint requires the support of the columns of $W$ to be disjoint as well, imposing $W$ to learn different features from the data set. Equivalently, it implies that there is at most a single non-zero entry in each row of $\mathcal{A}$ and $W$.

Looking back at the model $Y = AX \approx (AW)H = \mathcal{A}H$, this means that each row of $Y$ is approximated as an scaling of a single row of $H$ [33]. Because we will not enforce (5d) strictly, but use a regularization, our minimization problem is related to soft clustering, which can be effectively used to capture diurnal similarity patterns (temporal dependencies) present in the link count matrix $Y$.

*Estimation of $x_t$.* Another salient aspect of the proposed model is that the estimate of an OD flows, $\hat{x}_t$ at time $t$, using the observed link count $y_t$, is purely nonnegative. In fact, as explained in Section 1, $\hat{x}_t$ is estimated from the model

$$y_t \approx Ax_t = AWh_t = \mathcal{A}h_t,$$

and by taking $h_t = \mathcal{A}^\dagger y_t$. We have $\mathcal{A}^\dagger = \mathcal{A}^T$ since $\mathcal{A}^T \mathcal{A} = I_k$, so that $x_t = Wh_t = W\mathcal{A}^T y_t$ which is nonnegative since $A, \mathcal{A}, y_t \geq 0$. This desirable feature of the proposed method is due to the additional nonnegativity and orthogonality constraints, (5b) and (5d). To the best of the authors' knowledge, this aspect has been missing in current state-of-the-art methods that generally employ an additional step of setting negative entries to zero.

### 2.2. Temporal modeling using autoregression

The constraint (5c) defines $k$ independent autoregression models for $k$ timeseries of latent flows corresponding to the $k$ rows of the matrix $H$. For the $p$th timeseries of latent flows, the autoregression model approximates every element in that timeseries as a weighted sum of its previous elements, that is,

$$H(p, t) = \sum_{l \in \mathcal{L}} \Omega(p, l)H(p, t - l).$$

When computing the weighted sum for an element $H(p, t)$, not all elements preceding it are taken into account but only the elements $H(p, t - l)$ for $l \in \mathcal{L} = \{l : l \in \mathbb{N} \text{ and } l \ll T\}$, where $\mathcal{L}$ is the

6

lag set. In our model (5c), all $k$ autoregressive models share a common lag set $\mathcal{L}$, but there are $k$ weight vectors gathered in the matrix $\Omega$. The lag $\mathcal{L}$ have indices to indicate positive correlation among OD flows. In particular, the OD flows $\boldsymbol{x}_{t-l}$ and $\boldsymbol{x}_t$ are assumed to be correlated if $l \in \mathcal{L}$. Equivalently, this means that the OD flow vectors $\boldsymbol{x}_{t-l}$ and $\boldsymbol{x}_t$ are assumed to be correlated, since our model assumes $\boldsymbol{x}_t \approx W\boldsymbol{h}_t$ for all $t$.

Modeling temporal dependencies via (5c), that is, using $k$ autoregressive models sharing a common lag set $\mathcal{L}$, has several advantages over using a multivariate autoregressive model or over conventional graph embedding approach:

1. It requires only $k \times |\mathcal{L}|$ nonnegative weights in contrast to multivariate autoregressive models which require $k \times k \times |\mathcal{L}|$ nonnegative weights. Hence, it is less prone to overfitting and noise.

2. There is no restriction on defining indices of the lag set $\mathcal{L}$. In contrast, conventional graph embedding approaches (e.g. [24]) either use a lag set with few elements with short dependencies, or require a prior knowledge of temporal dependencies which is generally not available.

3. Unlike conventional graph embedding approaches (e.g. [24]), weights of $k$ autoregression models (that is, $\Omega$) can be learnt dynamically.

4. Like conventional graph embedding approaches such as [24], each of the $k$ autoregression models have an equivalent graph representation as explained in [26]; see Figure 1(a). Hence, the constraint (5c) can be incorporated into the data fitting term (5a) using exiting
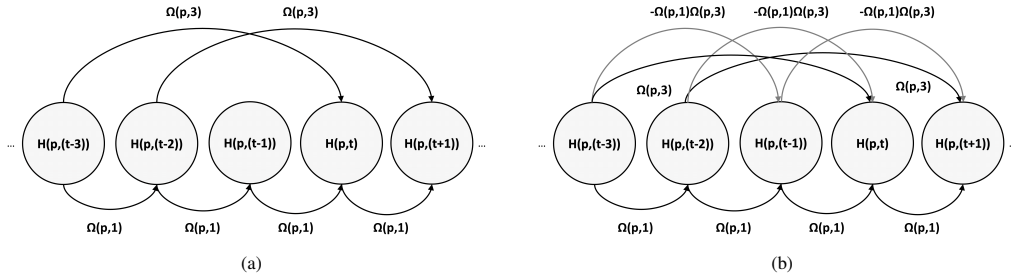


Figure 1: (a) A nonnegative weighted graph $\mathcal{G}$ defined by a typical graph embedding approach and associated to the $p$th row of $H$. (b) A weighted signed graph $\mathcal{G}^{AR}$ defined by the $p$th autoregressive model $H(p,t) = \sum_{l \in \mathcal{L}} \Omega(p,l)H(p,t-l)$ associated to the $p$th row of $H$. In both cases the lag set is supposed to be $\mathcal{L} = \{1, 3\}$. In (b) gray colored arrows indicates negatively weighted edges due to negative correlation between the two nodes. This figure is inspired from [26, Figures 2-3].

   techniques like Laplacian regularizer.

5. Unlike conventional graph embedding approaches such as [24], the graph associated with the $p$th autoregression model may contain negatively weighted edges to indicate negative correlation between two nodes of the graph; see Figure 1(b) for an illustration.

We refer the interested readers to [26] and the references for more details.

### 3. Training our traffic estimation model (5)

In this section, we first provide our regularized NMF model to tackle (5) in Section 3.1. In fact, in practice, because of noise and model misfit, it is not reasonable to strictly enforce the constraints (5c) and (5d) (in fact, the orthogonality constraint could even make the problem infeasible), and hence it makes more sense to only penalize the solutions that violate these constraints. Then we propose a fast gradient method to solve it in Section 3.2.

#### 3.1. Incorporation of constraints into the model

Let us consider the following regularized NMF-based network traffic estimation model, which replaces the constraints (5c) and (5d) in (5) with regularization terms:

$$\min_{W \in \mathbb{R}_+^{n \times k}, H \in \mathbb{R}_+^{k \times T}, \Omega \in \mathbb{R}_+^{k \times L}} F(W, H, \Omega), \tag{6}$$

$$\text{where } F(W, H, \Omega) = \|X - WH\|_F^2 + \lambda_h \sum_{p=1}^{k} \mathcal{T}(\omega_p, \mathcal{L}, h_p) + \lambda_{\mathcal{A}} \mathcal{I}(\mathcal{A}),$$

where $\mathcal{T}(\omega_p, \mathcal{L}, h_p)$ is the temporal regularizer to incorporate the constraint (5c) with $\omega_p$ (resp. $h_p$) the $p$th row of $\Omega$ (resp. $H$), and $\mathcal{I}(\mathcal{A})$ is the orthogonality regularizer to incorporate the constraint (5d); see below for more details. The positive numbers $\lambda_h$ and $\lambda_{\mathcal{A}}$ are the penalty parameters for the regularizers $\mathcal{T}(\omega_p, \mathcal{L}, h_p)$ and $\mathcal{I}(\mathcal{A})$ respectively and, $\omega_p$ and $h_p$ are the $p$th rows of $\Omega$ and $H$ respectively.

Let us briefly discuss the two regularizers.

*Orthogonality.* The orthogonality regularizer $\mathcal{I}(\mathcal{A})$ in (6) corresponds to the orthogonality constraint (5d) and is defined as:

$$\mathcal{I}(\mathcal{A}) = \|\mathcal{A}^\top \mathcal{A} - I\|_F^2.$$

This is standard least-squares penalty for orthogonality constraints; see for example [34].

*Autoregression.* We use the $p$th temporal regularizer $\mathcal{T}(\omega_p, \mathcal{L}, h_p)$ from [26], which is equivalent to the following Laplacian regularizer:

$$\mathcal{T}(\omega_p, \mathcal{L}, h_p) \quad = \quad \frac{1}{2} \sum_{t_1, t_2}^{T} S_p^{AR}(t_1, t_2) \|H(p, t_1) \quad - \quad H(p, t_2)\|^2 \quad + \quad \frac{1}{2} h_p D_p h_p^\top \tag{7}$$

where $S_p^{AR} \in \mathbb{R}^{T \times T}$ is the weight matrix of the graph $\mathcal{G}^{AR}$ and $D_p \in \mathbb{R}^{T \times T}$ is the diagonal matrix such that:

$$S_p^{AR}(t, t + d) = \begin{cases} \sum_{l \in \delta(d)} \sum_{L < t \le T} -\Omega(p, l)\Omega(p, (l - d)), & \text{if } \delta(d) \ne \phi, \\ 0, & \text{otherwise,} \end{cases}$$

where $\delta(d) = \{l \in \mathcal{L} \cup \{0\} : l - d \in \mathcal{L} \cup \{0\}\}$ and $L = \max(\mathcal{L})$, and

$$D_p(t, t) = \left( \sum_{l \in \mathcal{L} \cup \{0\}} \Omega(p, l) \right)\left( \sum_{l \in \mathcal{L} \cup \{0\}} \Omega(p, l)[L < t + l \le T] \right).$$

The graph Laplacian $L_p^{AR} \in \mathbb{R}^{T \times T}$ is defined as:

$$L_p^{AR}(t_1, t_2) = \begin{cases} \sum_{t_3=1}^{T} S_p^{AR}(t_1, t_3) & \text{if } t_1 = t_2, \\ -S_p^{AR}(t_1, t_2) & \text{otherwise.} \end{cases}$$

The first term in (7) is standard for Laplacian regularizers of conventional graph embedding approaches (see, e.g., [24]), whereas the second term is specific to the autoregression model [26] due to negative edges in the associated graph $\mathcal{G}^{AR}$ indicating negative correlation between nodes; see Figure 1b.

### 3.2. Fast gradient algorithm

In this section, we describe a training algorithm for the network traffic estimation model (6). The algorithm is iterative and is based on block coordinate descent method, as most NMF algorithms [35]. Such iterative algorithms require an initialization. The initial iterates are denoted $W^{(0)}$, $H^{(0)}$ and $\Omega^{(0)}$ which we compute as follows:

1. For the given training OD flow matrix $X$ and the rank of factorization $k$, the matrices $W$ and $H$ are initialized by $(W^{(0)}, H^{(0)}) = \text{NNSVD-LRC}(X, k)$ where NNSVD-LRC is an effective initialization for NMF based on the SVD [30].

2. The $p$th row $\omega_p$ of $\Omega$ is initialized as the projection of the unconstrained solution:

$$\omega_p^{(0)} = \mathcal{P}_+(h_p \mathcal{H}_p^{\dagger}),$$

where $\mathcal{P}_+(Z)$ is the projection onto the nonnegative orthant, that is,

$$\mathcal{P}_+(Z)_{i,j} = \begin{cases} Z(i, j), & \text{if } Z(i, j) \geq 0, \\ 0, & \text{otherwise,} \end{cases}$$

where $\mathcal{P}_+(Z)_{i,j}$ denotes the entry of $\mathcal{P}_+(Z)$ at position $(i, j)$, $\mathcal{H}_p^{\dagger}$ is the (left) pseudoinverse of $\mathcal{H}_p$ and $\mathcal{H}_p \in \mathbb{R}^{L \times T}$ is defined as:

$$\mathcal{H}_p(q, t) = \begin{cases} H(p, (t - q)), & \text{if } q \in \mathcal{L} \text{ and } L < t \leq T \\ 0, & \text{otherwise.} \end{cases}$$

Given the initial estimate $(W^{(0)}, H^{(0)}, \Omega^{(0)})$, Algorithm 1 further improves the solution via an iterative process. We use a standard strategy in NMF, that is, alternatively updating the block variables $(W, H, \Omega)$. For each block, we use a first-order accelerated gradient descent method with optimal convergence rate [36]; as done in [37] for the standard NMF model. In a nutshell, such methods take a gradient step from an extrapolated sequence.

Algorithm 1 describes the proposed training algorithm in detail. It takes as as input the initial values of three free parameters $W^{(0)}$, $H^{(0)}$ and $\Omega^{(0)}$, the training OD flows $X$, the routing matrix $A$, the rank of factorization $k$, the lag set $\mathcal{L}$, and the maximum number of iterations $n_{max}$. In each iteration Algorithm 1 alternatively optimizes $W$, $H$ and $\Omega$ using the accelerated gradient method as described in Algorithm 2 with help of Table 1. Note that we use a restarting mechanism in the fast gradient method which typically performs well in practice; see the discussion in [38].

Because accelerated gradient methods do not ensure the objective function to decrease at every iteration, and because we optimize alternatively $(W, H, \Omega)$, we embed Algorithm 2 with a restarting scheme: if the objective function increases, the algorithm abandon the extrapolated sequence and takes a standard gradient step, as done for example in [39]. This ensures the objective function will decrease.

9

---

**Algorithm 1** Solving the traffic matrix estimation model (6), MCST-NMF

---

**Input:** $X$ is an $n$-by-$T$ training OD flow matrix, $W^{(0)}, H^{(0)}$ and $\Omega^{(0)}$ are the initial values for $W, H$ and $\Omega$ respectively, $q^{max}$ is the maximum number of iterations, and $0 \leq \delta \ll 1$ is a threshold for early stopping when the differences between two consecutive errors is small enough.

**Output:** Final values for $W$, $H$ and $\Omega$ and $\epsilon^{min}$ solving (6)

1: Initialize: $W = W^{(0)}$; $H = H^{(0)}$; $\Omega = \Omega^{(0)}$; $e(0) = \|X - WH\|_F^2$;
   $\quad\quad \epsilon_{min} = \delta e(0)$; $\epsilon = \epsilon^{min}$; $q = 1$;
2: **repeat**
3: $\quad W = \text{FastGradientUpdate}(X, W, H, \Omega)$;
4: $\quad H = \text{FastGradientUpdate}(X, W, H, \Omega)$;
5: $\quad \Omega = \text{FastGradientUpdate}(X, W, H, \Omega)$;
6: $\quad e(q) = \|X - WH\|_F^2$; $\epsilon = e(q-1) - e(q)$; $q = q + 1$;
7: **until** $q \leq q^{max}$ and ($\epsilon < 0$ or $\epsilon \geq \epsilon^{min}$)

---

**Algorithm 2** FastGradientUpdate$(X, W, H, \Omega)$

---

**Input:** $X$ is an $n$-by-$T$ training OD flow matrix, $W$, $H$, and $\Omega$ are the three free parameters, $q_B^{max}$ the maximum iterations for loop, and $1 \gg \delta_B > 0$ the minimum differences between two consecutive errors w.r.t. initial error for early stopping before $q_B^{max}$ iterations respectively.

**Output:** Updated $B$ by solving (6) for $B$ (where $B$ can be $W$, $H$ or $\Omega$)

1: Initialize: $C = B$; $B^{(curr)} = B$; $\alpha_B = \alpha_B^{(prev)} = 1$; $q_B = 1$;
2: Compute $e_B^{(curr)}$ as described in Table 1
3: $e_B^{(prev)} = e_B^{(curr)}$; $\epsilon_B^{min} = \delta_B e_B^{(prev)}$;
4: **repeat**
5: $\quad \alpha_B = \frac{1 + \sqrt{4\alpha_B^2 + 1}}{2}$;
6: $\quad$ Compute $L_B$ and $\nabla_B F(W, H, \Omega)$ as described in Table 1.
7: $\quad B = \mathcal{P}_+\left(C - \frac{1}{L_B} \nabla_B F(W, H, \Omega)\right)$;
8: $\quad C = B + \frac{\alpha_B^{(prev)} - 1}{\alpha_B}\left(B - B^{(curr)}\right)$;
9: $\quad$ Compute $e_B^{(curr)}$ as described in Table 1
10: $\quad \epsilon_B = e_B^{(prev)} - e_B^{(curr)}$;
11: $\quad$ **if** $\epsilon_B < 0$ **then**
12: $\quad\quad C = B$; $\alpha_B = 1$; $e_B^{(curr)} = e_B^{(prev)}$;
13: $\quad$ **end if**
14: $\quad \alpha_B^{(prev)} = \alpha_B$; $e_B^{(prev)} = e_B^{(curr)}$; $q_B = q_B + 1$;
15: **until** $q_B \leq q_B^{max}$ and ($\epsilon_B < 0$ or $\epsilon_B \geq \epsilon_B^{min}$)

---

*Choice of the penalty parameters.* It is not an easy task to tune the two penalty parameters $\lambda_h$ and $\lambda_{\mathcal{A}}$ for the regularization terms in (6). Given the initial iterate $(W^{(0)}, H^{(0)}, \Omega^{(0)})$, we will use

$$\lambda_h = \beta_h \frac{\|X - W^{(0)} H^{(0)}\|_F^2}{\sum_{p=1}^k \|h_p - \omega_p^{(0)} \mathcal{H}_p^{(0)}\|_2^2}, \tag{8a}$$

$$\lambda_{\mathcal{A}} = \beta_{\mathcal{A}} \frac{\|X - W^{(0)} H^{(0)}\|_F^2}{\|(\mathcal{A}^{(0)})^\top \mathcal{A}^{(0)} - I\|_2^2}, \tag{8b}$$

Table 1: Computations needed for Algorithm 2, FastGradientUpdate$(X, W, H, \Omega)$, depending on the updated variable, $B$ is either $W$, $H$ or $\Omega$.

| variable | quantity | formula |
|---|---|---|
| $W$ | $L_B$ | $\|W^\top W\|_2$ |
| | $\nabla_B F(W, H, \Omega)$ | $2(WHH^\top - XH^\top) + 4\lambda_{\mathcal{A}}(A^\top \mathcal{A})(\mathcal{A}^\top \mathcal{A} - I_k)$ |
| | $e_B^{(curr)}$ | $\|X - WH\|_F^2$ |
| $H$ | $L_B$ | $\|HH^\top\|_2 + \sum_{p=1}^k \lambda_h(\|L_p^{AR}\|_2 + \|D_p\|_2)$ |
| | $\nabla_B F(W, H, \Omega)$ | $2(W^\top WH - W^\top X) + \lambda_h H(p, :)\big((L_p^{AR})^\top + L_p^{AR}\big)$ |
| | $e_B^{(curr)}$ | $\|X - WH\|_F^2$ |
| $\Omega$ | $L_{B(p,:)}$ | $\|\mathcal{H}_p \mathcal{H}_p^\top\|_2, \qquad$ for $p = 1, 2, \ldots, k.$ |
| | $\nabla_{B(p,:)} F(W, H, \Omega)$ | $2\big(\Omega(p, :)\mathcal{H}_p \mathcal{H}_p^\top - H(p, :)\mathcal{H}_p^\top\big)$ |
| | $e_B^{(curr)}$ | $\sum_{p=1}^k \|H(p, :) - \Omega(p, :)\mathcal{H}_p\|_2^2$ |

where $\beta_h, \beta_{\mathcal{A}} \in (0, 1]$. This choice allows to balance the importance of the penalty terms compared to the data fitting term, at initialization. Like other existing regularized models, e.g., [25], the tuning of the penalty parameters, $\beta_h$ and $\beta_{\mathcal{A}}$ in our model, is a difficult task and typically problem dependent (nature of the given network, noise level, etc.). In practice, a useful way to tune parameters is to use cross validation (see, e.g., [40]), that is, use training sets to train a model with different parameter values, and then select the values of the parameters that lead to the best results on the test sets. For our problem, a value of $\beta_h$ and $\beta_{\mathcal{A}}$ around 10-20% typically works well in practice.

*Handling the missing entries.* It is usual that the traffic matrix $X$ used for training a given traffic estimation model contains missing entries, even in the presence of a good measurement system [7]. To the best of our knowledge, all the existing state-of-the-art algorithms assume that the training data is full and complete. To address the issue of an incomplete training dataset, we propose two approaches. The first approach preprocesses the given training dataset to fill in the missing entries using a weighted NMF model trained using a fast gradient method [41]. The second approach modifies Algorithm 1 using the expectation maximization strategy proposed in [42].

Denoting $M \in \{0, 1\}^{n \times T}$ the binary mask matrix so that $M(i, j) = 1$ if and only if $X(i, j)$ is observed, the second approach for handling the missing entries modifies Algorithm 1 to use the data matrix $X^{(q)}$ instead of $X$ at iteration $q$, where $X^{(q)}$ is defined as

$$X^{(q)} = M \circ X + (\mathbb{1}_{n \times T} - M) \circ (W^{(q-1)} H^{(q-1)})$$

where $\circ$ is the element-wise matrix multiplication and $\mathbb{1}_{n \times T}$ is the $n \times T$ all one matrix.

### 3.3. Estimation of future traffic flows

It is a common practice to improve the solution obtained by the traffic estimation model using an Expectation Maximization (EM) iterative algorithm or iterative proportional fitting (IPF) algorithm; see for example [9, 14, 15]. In this paper, we use the EM approach from [43] combined

11

---

**Algorithm 3** Estimating OD flows

---

**Input:** $A$ an $m$-by-$n$ routing matrix, $W$ the $n$-by-$k$ matrix, the new observed link flow $\boldsymbol{y}_t$, $q^{max}$, and $r^{max}$ the maximum number of iterations for fast gradient descent and expectation maximization iteration steps respectively, $1 \gg \delta_{gd}, \delta_{emi} > 0$ are the minimum difference between two consecutive errors for early stopping before $q^{max}$ and $r^{max}$ iterations of fast gradient descent expectation maximization iteration steps, respectively.

**Output:** Final estimated OD flow $\hat{\boldsymbol{x}}_t$ for timestamp $t$.

1: Initialization for fast gradient steps:
2: $\mathcal{A} = AW$; $\boldsymbol{h}_t = \mathcal{A}^\top \boldsymbol{y}_t$; $\boldsymbol{v}_t = \boldsymbol{h}_t$; $\boldsymbol{h}_t^{(curr)} = \boldsymbol{h}_t$;
3: $\alpha_{gd}^{(prev)} = 1$; $\eta = \frac{1}{\|\mathcal{A}^\top \mathcal{A}\|_2}$; $e_{gd}^{(prev)} = \|\boldsymbol{y}_t - A\boldsymbol{h}_t\|_2^2$;
4: $\epsilon_{gd}^{min} = \delta_{gd}\|y_t\|_2^2$; $q = 1$;
5: **repeat**
6: $\quad \alpha_{gd} = \frac{1+\sqrt{4\alpha_{gd}^2+1}}{2}$;
7: $\quad \boldsymbol{h}_t = \mathcal{P}_+\left(\boldsymbol{v}_t - \eta\nabla_{\boldsymbol{h}_t}G(\boldsymbol{y}_t, \mathcal{A}, \boldsymbol{v}_t)\right)$;
8: $\quad \boldsymbol{v}_t = \boldsymbol{h}_t + \frac{\alpha_{gd}^{(prev)}-1}{\alpha_{gd}}\left(\boldsymbol{h}_t - \boldsymbol{h}_t^{(curr)}\right)$;
9: $\quad e_{gd}^{(curr)} = \|\boldsymbol{y}_t - A\boldsymbol{v}_t\|_2^2$; $\epsilon_{gd} = e_{gd}^{(prev)} - e_{gd}^{(curr)}$;
10: $\quad$ **if** $\epsilon_{gd} < 0$ **then**
11: $\quad\quad \boldsymbol{v}_t = \boldsymbol{h}_t$; $\alpha_{gd} = 1$; $e_{gd}^{(curr)} = e_{gd}^{(prev)}$;
12: $\quad$ **end if**
13: $\quad \alpha_{gd}^{(prev)} = \alpha_{gd}$; $e_{gd}^{(prev)} = e_{gd}^{(curr)}$; $q = q + 1$;
14: **until** $q \le q^{max}$ and $(\epsilon_{gd} < 0$ or $\epsilon_{gd} \ge \epsilon_{gd}^{min})$
15: Initialization for expectation maximization steps:
16: $\boldsymbol{x}^{(curr)} = W\boldsymbol{h}_t$; $\boldsymbol{y} = \boldsymbol{y}_t$;
17: $e_{emi}(0) = 0$; $\epsilon = \epsilon_{min}$; $\epsilon_{emi}^{min} = \delta_{emi}\|\boldsymbol{x}^{(curr)}\|_2^2$; $r = 1$;
18: **repeat**
19: $\quad$ **for** $j = 1 : 1 : n$ **do**
20: $\quad\quad \boldsymbol{x}(j) = \frac{x(j)^{(curr)}}{\sum_{i=1}^m A(i,j)}\sum_{i=1}^m \frac{A(i,j)y(i)}{\sum_{k=1}^n A(i,k)x(k)^{(curr)}}$
21: $\quad$ **end for**
22: $\quad \epsilon_{emi} = \|\boldsymbol{x} - \boldsymbol{x}^{(curr)}\|_2^2$; $\boldsymbol{x}^{(curr)} = \boldsymbol{x}$; $r = r + 1$;
23: **until** $r \le r^{max}$ and $(\epsilon_{emi} \ge \epsilon_{emi}^{min})$
24: $\hat{\boldsymbol{x}}_t = \boldsymbol{x}$;

---

with a fast gradient method, which is empirically found to be more effective than applying the EM algorithm alone. Algorithm 3 provides the exact details of our traffic estimation procedure.

Once the proposed model (6) has been trained using Algorithm 1 over a given network, let us explain how the new unobserved OD flows $\boldsymbol{x}_t$ can be estimated as $\hat{\boldsymbol{x}}_t$ using the estimated parameters $W$, the routing matrix $A$, and the observed link flow $\boldsymbol{y}_t$; see Algorithm 3. The algorithm first estimates the latent flow $h_t^{(0)}$ as $\mathcal{A}^T\boldsymbol{y}_t$; see Section 2.1. It is then refined using a few steps of projected fast gradient descent applied on the following minimization problem:

$$\min_{\boldsymbol{h}_t \in \mathbb{R}_+^k} G(\boldsymbol{y}_t, \mathcal{A}, \boldsymbol{h}_t) = \min_{\boldsymbol{h}_t \in \mathbb{R}_+^k} \|\boldsymbol{y}_t - \mathcal{A}\boldsymbol{h}_t\|_2^2.$$

Using the final estimated latent flow $\boldsymbol{h}_t^{(q_f)}$ after $q_f$th iteration of the projected fast gradient

12

method, initialize $\boldsymbol{x}^{(0)} = W\boldsymbol{h}_t^{(q_f)}$ and $\boldsymbol{y} = \boldsymbol{y}_t$. The $r$th iteration of EM for solving (1) is as follows [43, 44]:

$$\boldsymbol{x}(j)^{(r)} = \frac{\boldsymbol{x}(j)^{(r-1)}}{\sum_{i=1}^{m} A(i, j)} \sum_{i=1}^{m} \frac{A(i, j)\boldsymbol{y}(i)}{\sum_{k=1}^{n} A(i, k)\boldsymbol{x}(k)^{(r-1)}}, \quad \forall\ j = 1, 2, \ldots, n,$$

where $\boldsymbol{x}(j)^{(r-1)}$ and $\boldsymbol{x}(k)^{(r-1)}$ are the $j$th and $k$th OD flows of $\boldsymbol{x}^{(r-1)}$ respectively, $\boldsymbol{y}(i)$ is the $i$th link flow of $\boldsymbol{y}$ and $A(i, j)$ is the $(i, j)$th element of the binary routing matrix $A$. The output of the EM algorithm after $r_f$th iteration will be the final estimated OD flow $\hat{\boldsymbol{x}}_t = \boldsymbol{x}^{(r_f)}$ and be the output of Algorithm 3.

## 4. Experimental Results

In this section, we evaluate the performance of our proposed traffic estimation model which we refer to as MCST-NMF (if the training dataset has missing entries, we use W-NeNMF [41] by default to fill in the missing entries), otherwise as MCST-NMC when missing entries are dealt with using EM. Moreover, we consider the following three state-of-the-art methods for the performance comparison a) PCA [9], b) MNETME [15], and c) CS-DME [13]. We consider these state-of-the-art because they are all based on linear dimensionality reduction, like MCST-NMF; see Section 1.

*Performance metrics.* Given the true $X$ the true OD flow matrix of a test data and the corresponding estimated OD flow matrix $\widehat{X}$ by an algorithm, we use the two performance metrics SRE and TRE. They are row and column vectors of dimension $n$ and $T$ respectively. The $i$th and $t$th elements of SRE and TRE vectors are defined as follows:

$$SRE(i) = \frac{\sqrt{\sum_{t=1}^{T} \left(\widehat{X}_{i,t} - X_{i,t}\right)^2}}{\sqrt{\sum_{t=1}^{T} X_{i,t}^2}}, \tag{9}$$

$$TRE(t) = \frac{\sqrt{\sum_{i=1}^{n} \left(\widehat{X}_{i,t} - X_{i,t}\right)^2}}{\sqrt{\sum_{i=1}^{n} X_{i,t}^2}}. \tag{10}$$

The scalar $SRE(i)$ is the normalized mean squared error of $i$th OD flow among all $T$ test timestamps, whereas $TRE(t)$ is the normalized mean squared error of the $t$th test timestamp among all $n$ OD flows.

These two performance metrics, SRE and TRE, are widely used in the literature to measure the performance a given network traffic estimation model. $SRE(i)$ primarily describes how well a given model is able to estimate a given OD flow $i$ over a given time period $T$. Hence, it provides us the insight about the performance of the given model w.r.t. a given OD flow. In contrast, $TRE(t)$ describes the efficiency of a given model to estimate all OD flows in the given network at a particular timestamp $t$. Hence, it provides us with the insight about the performance of the given model w.r.t. a given timestamp.

Tables 2 and 4 report the statistical properties of these two vectors, that is, of the SRE and TRE on the test dataset.

*Datasets.* The performance of different models is evaluated over two publicly available and widely used Internet traffic datasets, namely Internet2 [31] and GÉANT [32]. Internet2 (a.k.a. Abilene) is the high-speed backbone network of US. It consists of 12 routers, 144 OD flows and, 15 and 12 bidirectional internal and external links respectively. The topology of Internet2 network is depicted in Figure 2-(a). The dataset [31] provides the measurements of 24 weeks of OD flows and the routing matrix of the network. The traffic is measured at an interval of 5 min and reported in the unit of 100 bytes. GÉANT is a research and educational network for Europe which is almost twice as large as Internet2. Figure 2-(b) illustrates the topology of the network. It has 23 routers, $23 \times 23 = 529$ OD flows and, 37 and 23 bidirectional internal and external links respectively. GÉANT [32] is a collection of 4 months of OD flows measurements and the routing information of the network. The corresponding link flow matrix $Y$ can be obtained using (2). Traffic measurements are in unit of kpbs and were performed at an interval of 15 minutes. Note that, in both cases, the link flow matrix $Y$ can be calculated using (2).
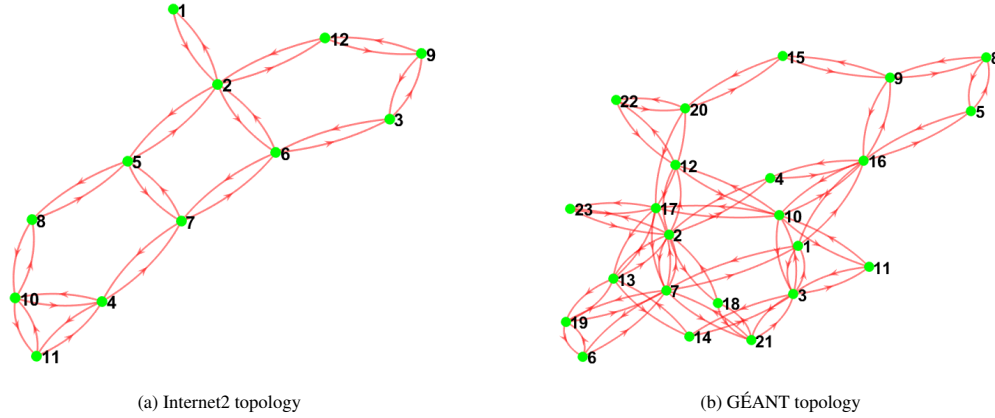


(a) Internet2 topology          (b) GÉANT topology

Figure 2: (a) Internet2 topology. (b) GÉANT topology. Routers are numbered green coloured circles. Links are red coloured arrows. Only internal links are shown.

*Parameter setting for experiments.* For experiments on the two datasets Internet2 and GÉANT, we use the following specifications:

i) For Internet2, the initial 11 days time slots ($11 * 24 * 12 = 3168$ time slots) of traffic data are used during experiments. The initial 7 days ($7 * 24 * 12 = 2016$ time slots) of traffic data in this time slot are used for training, and the subsequent 4 days ($4 * 24 * 12 = 1152$ time slots) traffic data are used for testing. For GÉANT, the three weeks of traffic data ($21 * 24 * 4 = 2016$ time slots) are used during experiments. The initial two weeks of traffic data ($14 * 24 * 4 = 1344$ time slots) are used for training of models, and the subsequent one week of traffic data ($7 * 24 * 4 = 672$ time slots) is used for testing.

ii) The factorization rank $k$ is set to 20 for all models.

iii) The models other than the proposed model are run with their default parameters.

iv) During training MCST-NMF, the penalty parameters $\lambda_h$ and $\lambda_{\mathcal{A}}$ for regularization terms are obtained using (8) by setting $\beta_h = \beta_{\mathcal{A}}$. For Internet2, $\beta_h = \beta_{\mathcal{A}} = 0.2$. For GÉANT, $\beta_h = 0.1$ and $\beta_A = 0.1$ respectively.

14

v) During training MCST-NMF, the maximum number of iterations $q^{max}$ is set to 50, while $q_W^{max}$, $q_H^{max}$ and $q_\Omega^{max}$ are set to 10. The parameter $\delta$ is set to $10^{-9}$, $\delta_W$ and $\delta_H$ are set to $10^{-3}$, and $\delta_\Omega$ is set to $10^{-5}$.

vi) For testing MCST-NMF, Algorithm 3 parameters are set as follows: a) $\delta_{gd}$ and $\delta_{emi}$ are set to $10^{-3}$ and $10^{-9}$, respectively, and b) $q^{max}$ and $r^{max}$ are set to 200.

*Choice of the lag set $\mathcal{L}$.* For training MCST-NMF using Algorithm 1 requires a lag set. We strategically choose the lag set during experiments by keeping the general temporal behaviour of backbone computer network in our mind i.e. a) network traffic behaviour persist for short period of time, b) network traffic is likely to repeat its behaviour on a hourly basis, c) network traffic may be similar after 8 hours at the beginning and end of working hours and d) network traffic exhibit diurnal pattern. Based on the described strategy, the chosen lag sets are $\mathcal{L} = \{1, 2, 3, 12, 24, 96, 102, 108, 288\}$ and $\mathcal{L} = \{1, 4, 8, 32, 34, 36, 96\}$ for Internet2 and GÉANT respectively.

All tests are preformed using MATLAB® R2018b (Student License) under Windows 10® environment on a laptop Intel® CORE™ i5-3YY6U768Y0M CPU @2.60GHz 4GB RAM. The code of our proposed methods is available from `https://github.com/5y3datif/MCST-NMF`.

### 4.1. Experiments on Internet2

Table 2: Statistical properties of the traffic estimation errors, SRE (9) and TRE (10), on the test set for the Internet2 data set. The lowest (best) values are highlighted in bold.

(a) Statistical properties of the SRE.

|  | MCST-NMC | MCST-NMF | MNETME | CS-DME | PCA |
|---|---|---|---|---|---|
| minimum | 0.07 | **0.05** | 0.07 | 0.81 | 0.07 |
| maximum | 1.36 | 1.34 | 11.86 | **1.03** | 2.37 |
| mean | 0.40 | **0.39** | 0.69 | 0.92 | 0.51 |
| median | **0.36** | **0.36** | 0.43 | 0.92 | 0.43 |
| standard deviation | 0.25 | 0.25 | 1.23 | **0.03** | 0.35 |

(b) Statistical properties of the TRE.

|  | MCST-NMC | MCST-NMF | MNETME | CS-DME | PCA |
|---|---|---|---|---|---|
| minimum | **0.05** | 0.06 | 0.13 | 0.80 | 0.08 |
| maximum | **0.30** | 0.31 | 0.37 | 0.97 | 0.54 |
| mean | **0.18** | **0.18** | 0.26 | 0.91 | 0.30 |
| median | **0.18** | 0.19 | 0.26 | 0.91 | 0.32 |
| standard deviation | 0.04 | 0.05 | 0.04 | **0.01** | 0.08 |

Let us examine the results on the Internet2 dataset to assess the performance of our two proposed methods compared to the state of the art. Table 2 presents the statistical properties of SRE and TRE of the different methods.

15

SRE of our proposed methods, MCST-NMF and MCST-NMC, is the best in terms of the minimum, mean, median and standard deviation with values 0.05, 0.39, 0.36 and 0.25, respectively. The closest competitor in terms of the minimum, mean and median of SRE has values larger by 40%, 31% and 19%, respectively. In terms of the maximum value of SRE, our methods are the second best (1.34 and 1.36). Similarly, TRE of MCST-NMC and MCST-NMF have the lowest values w.r.t. all four criteria, namely the minimum (0.05 and0.06), maximum (0.30 and 0.31), mean (0.18 and 0.18) and median (0.18 and 0.19). Its nearest competitor with respect to corresponding criteria has values larger by 33%, 19%, 44% and 44%, respectively.

To better visualize the differences, Figure 3 displays the cumulative distribution functions (CDF) of the SRE and TRE of the different methods. We observe that 90% of OD flows estimated by our proposed methods (MCST-NMC and MCST-NMF) have SRE of at most 0.74. In contrast, the nearest competitor, PCA, has SRE of at most 0.90 for 90% of OD flows that is 21.62% higher comparing to that of MCST-NMF. Results for the TRE for MCST-NMC and MCST-NMF are even more striking. All the models show low TRE but 90% of OD flows estimated by our proposed methods (MCST-NMC and MCST-NMF) has TRE of at most 0.238 and 0.242, respectively, which is 6.7% and 5% smaller than the second best model, namely CS-DME (with value 0.254) as shown on Figure 3b.
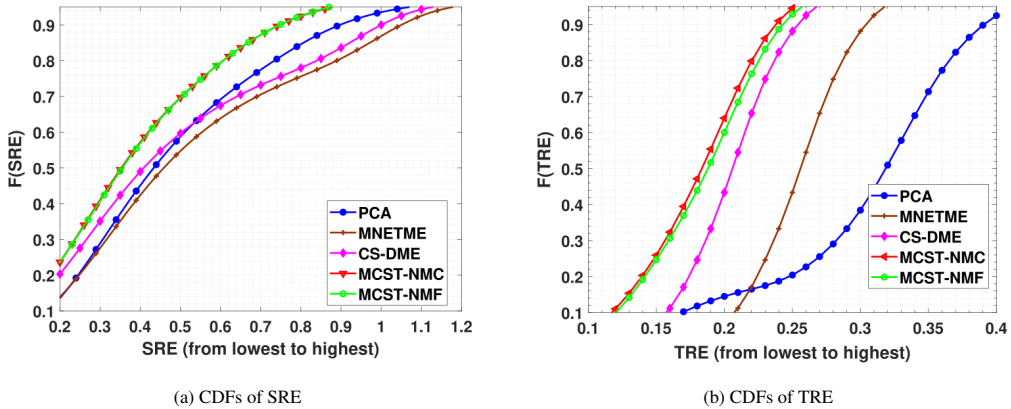


(a) CDFs of SRE             (b) CDFs of TRE

Figure 3: Cumulative distributed function (CDF) of SRE (left) and TRE (right) resulting from comparing traffic estimation methods over Internet2.

### 4.1.1. Impact of the regularization parameters

This section analyzes the impact of the regularization parameters, $\beta_h$ and $\beta_{\mathcal{A}}$, on the performance of the proposed methods when tested over the Internet2 dataset. To do so, we consider four cases by setting $\beta_h$ and $\beta_{\mathcal{A}}$ to their default value 0.5 or to 0, with a total of four cases. All the remaining parameters are kept as described in the paragraph Parameter setting for experiments during all the conducted experiments. The obtained results are reported in Table 3.

We observe that our two methods have the best performance when both the regularization terms are used. Using only one regularization term or none of the two regularization terms deteriorate SRE and TRE upto 24% and 7% respectively. This illustrates the effectiveness of the two regularization terms. Note that the case $\beta_h = 0$ represent the case where the lag set $\mathcal{L} = \emptyset$ which has a significant impact on the TRE; see Table 3.

Table 3: Impact of the regularization parameters with respect to SRE and TRE for the Internet2 dataset. Lowest values are in bold.

| | SRE | | TRE | |
| | MCST-NMF | MCST-NMC | MCST-NMF | MCST-NMC |
|---|---|---|---|---|
| $\beta_h = 0.5, \beta_{\mathcal{A}} = 0.5$ | **0.73** | **0.73** | 0.187 | **0.186** |
| $\beta_h = 0.5, \beta_{\mathcal{A}} = 0.0$ | 0.76 | 0.76 | 0.232 | 0.230 |
| $\beta_h = 0.0, \beta_{\mathcal{A}} = 0.5$ | 0.76 | 0.77 | 0.227 | 0.216 |
| $\beta_h = 0.0, \beta_{\mathcal{A}} = 0.0$ | 0.78 | 0.78 | 0.216 | 0.216 |

*4.2. Experiment on GÉANT*

We now perform the same results as for Internet2 on the GÉANT data set. Table 4 presents the values of SRE and TRE for the different methods.

Table 4: Statistical properties of SRE and TRE for PCA, MNETME, CS-DME and MCST-NMF over GÉANT test data. Lowest values are in bold.

(a) Statistical properties of SRE

| | MCST-NMC | MCST-NMF | MNETME | CS-DME | PCA |
|---|---|---|---|---|---|
| minimum | **0.00** | **0.00** | **0.00** | 0.12 | 0.02 |
| maximum | **23** | 46 | 59 | 30 | 37 |
| mean | 5.71 | 14.21 | 14.87 | **2.53** | 10.47 |
| median | **0.82** | 0.86 | 0.87 | 1.27 | 0.91 |
| standard deviation | 6.84 | 22.34 | 14.76 | **3.75** | 86.99 |

(b) Statistical properties of TRE

| | MCST-NMC | MCST-NMF | MNETME | CS-DME | PCA |
|---|---|---|---|---|---|
| minimum | **0.05** | **0.05** | 0.06 | 0.41 | 0.06 |
| maximum | **0.26** | **0.26** | 0.60 | 1.13 | 0.38 |
| mean | **0.09** | **0.09** | 0.13 | 0.64 | 0.12 |
| median | **0.08** | **0.08** | 0.10 | 0.61 | 0.11 |
| standard deviation | **0.03** | **0.03** | 0.10 | 0.12 | 0.05 |

The SRE of our proposed methods is the best in terms of the median. Moreover, the matrix completion variant of our proposed model, MCST-NMC, is the best in terms of all criteria except except for the mean and standard deviation. The closest competitors of MCST-NMC in terms of the minimum, maximum and median of SRE has values larger by 0%, 30% and 4.87%, respectively. The TREs of our proposed methods are the same. The closest competitor in terms of the minimum, maximum, mean, median and standard deviation has value larger by 20%, 46%, 33%, 25% and 66%, respectively. To further illustrate the differences, Figure 4 shows the cumu-

lative distribution functions (CDF) of the different methods. We observe that 80% of OD flows estimated by our proposed methods, MCST-NMC and MCST-NMF, have SRE of value at most 1.02 and 1.04, respectively. The second best performing model (MNETME) has SRE of at most 1.1 for 80% of OD flows. This value is larger by 8% compared to MCST-NMC; see Figure 4a. Results of TRE for MCST-NMC and MCST-NMF are very similar. All the models show low TRE but 90% of OD flows estimated by MCST-NMC and MCST-NMF have TRE of at most 0.12 whereas the second best model, PCA, has TRE of at most 0.15 for 90% of OD flows. This value is larger by 25%; see Figure 4b.



(a) CDFs of SRE

(b) CDFs of TRE

Figure 4: Cumulative distributed function (CDF) of SRE (left) and TRE (right) resulting from comparing traffic estimation models over GÉANT test data. (On the right plot, the TRE of the two proposed variants, MCST-NMF and MCST-NMC, overlap.)

### 4.2.1. Impact of the regularization parameters

This section analyzes the impact of the regularization parameters, $\beta_h$ and $\beta_{\mathcal{A}}$ on the performance of proposed solutions when tested over GÉANT dataset, exactly as for the Internet2 dataset. Table 5 reports the values of the SRE and TRE for the different values of the regularization parameters.

Table 5: Impact of the regularization parameters with respect to SRE and TRE over GÉANT. Lowest values are in bold.

| cases | SRE | | TRE | |
|---|---|---|---|---|
| | MCST-NMF | MCST-NMC | MCST-NMF | MCST-NMC |
| $\beta_h = 0.5, \beta_{\mathcal{A}} = 0.5$ | 1.20 | **1.12** | 0.124 | **0.121** |
| $\beta_h = 0.5, \beta_{\mathcal{A}} = 0.0$ | 1.27 | 1.21 | 0.125 | 0.123 |
| $\beta_h = 0.0, \beta_{\mathcal{A}} = 0.5$ | 4.68 | 1.70 | 0.124 | **0.121** |
| $\beta_h = 0.0, \beta_{\mathcal{A}} = 0.0$ | 4.25 | 1.80 | 0.124 | **0.121** |

We observe that both proposed methods exhibit the best performance when the regularization terms are used. This is particularly true for the SRE, with a significant deterioration for both methods; namely from 1.20 to 4.25 for MCST-NMF, and from 1.12 to 1.80 for MCST-NMC.

Moreover, we observe that the use of the lag set is particularly important to have low SRE values, as using $\beta_h = 0$ significantly impacts the SRE.

## 5. Conclusion

In this paper, we have proposed an NMF-based approach to tackle the network traffic flow estimation problem. To the best of our knowledge, it is the first time NMF is used for this specific task. A notable shortcoming of previously explored dimensionality-reduction approaches (e.g., based on PCA) was to solve these problems by ignoring the nonnegativity constraints. Moreover, our approach uses two regularizers, namely orthogonality to better cluster the data, and autoregression to take the temporal correlations into account, which further improves its performance.

We proposed their two different variants of our model: (1) MCST-NMF that is most suitable when the training dataset does not have missing entries, and (2) MCST-NMC that is designed specifically to handle missing entries. We have shown on two real-world data sets, namely GÉANT and Internet2, that our methods outperform existing techniques based on linear dimensionality reduction.

*Further Works.* Recently, deep matrix factorization models have emerged; in particular deep NMF has been shown to be able to capture several layers of meaningful features; see, e.g., [45, 46, 47]. It would be interesting to adapt our NMF based traffic matrix estimation model into a deep NMF model. Other directions of research include the use of other NMF models to perform the traffic flow estimation, such as [48, 49], or to improve prediction using data coming from other sources than the traffic flows and then use multi-view techniques such as [21, 50].

## Acknowledgement

## References

[1] K. Xie, L. Wang, X. Wang, G. Xie, J. Wen, G. Zhang, J. Cao, D. Zhang, Accurate Recovery of Internet Traffic Data: A Sequential Tensor Completion Approach, IEEE/ACM Trans. Networking 26 (2) (2018) 793–806. `doi:10.1109/TNET.2018.2797094`.
URL `http://ieeexplore.ieee.org/document/8301588/`

[2] F. Xiao, L. Chen, H. Zhu, R. Hong, R. Wang, Anomaly-Tolerant Network Traffic Estimation via Noise-Immune Temporal Matrix Completion Model, IEEE Journal on Selected Areas in Communications 37 (6) (2019) 1192–1204. `doi:10.1109/JSAC.2019.2904347`.
URL `https://ieeexplore.ieee.org/document/8664591/`

[3] C. Hu, S. Wang, J. Tian, B. Liu, Y. Cheng, Y. Chen, Accurate and Efficient Traffic Monitoring Using Adaptive Non-Linear Sampling Method, in: IEEE INFOCOM 2008 - IEEE Conference on Computer Communications, Phoenix, AZ, USA, 2008, pp. 26–30. `doi:10.1109/INFOCOM.2008.14`.
URL `http://ieeexplore.ieee.org/document/4509609/`

[4] P. Tune, D. Veitch, Fisher Information in Flow Size Distribution Estimation, IEEE Trans. Inform. Theory 57 (10) (2011) 7011–7035. `doi:10.1109/TIT.2011.2165150`.
URL `http://ieeexplore.ieee.org/document/6034747/`

[5] S. Pan, Y. Zhou, Z. Zhang, S. Yang, F. Qian, G. Hu, Identify Congested Links with Network Tomography Under Multipath Routing, Journal of Network and Systems Management 27 (2) (2019) 409–429. `doi:10.1007/s10922-018-9471-2`.
URL `http://link.springer.com/10.1007/s10922-018-9471-2`

[6] X. Li, K. Xie, X. Wang, G. Xie, J. Wen, G. Zhang, Z. Qin, Online Internet Anomaly Detection With High Accuracy: A Fast Tensor Factorization Solution, in: IEEE INFOCOM 2019 - IEEE Conference on Computer Communications, Paris, France, 2019, pp. 1900–1908. `doi:10.1109/INFOCOM.2019.8737562`.
URL `https://ieeexplore.ieee.org/document/8737562/`

[7] P. Tune, M. Roughan, H. Haddadi, O. Bonaventure, Internet traffic matrices: A primer, Recent Advances in Networking 1 (2013) 1–56.

[8] X. Wu, K. Yu, X. Wang, On the growth of Internet application flows: A complex network perspective, in: IEEE INFOCOM 2011 - IEEE Conference on Computer Communications, Shanghai, China, 2011, pp. 2096–2104. `doi:10.1109/INFOCOM.2011.5935019`.
URL `http://ieeexplore.ieee.org/document/5935019/`

[9] A. Soule, A. Lakhina, N. Taft, K. Papagiannaki, K. Salamatian, A. Nucci, M. Crovella, C. Diot, Traffic matrices: balancing measurements, inference and modeling, in: Proceedings of the 2005 ACM SIGMETRICS international conference on Measurement and modeling of computer systems - SIGMETRICS '05, ACM Press, Banff, Alberta, Canada, 2005, p. 362. `doi:10.1145/1064212.1064259`.
URL `http://portal.acm.org/citation.cfm?doid=1064212.1064259`

[10] G. H. Golub, C. F. Van Loan, Matrix computations, fourth edition Edition, Johns Hopkins studies in the mathematical sciences, The Johns Hopkins University Press, Baltimore, 2013, oCLC: ocn824733531.

[11] A. Kumar, V. V. Saradhi, T. Venkatesh, Compressive Sensing of Internet Traffic Matrices using CUR Decomposition, in: Proceedings of the 19th International Conference on Distributed Computing and Networking - ICDCN '18, ACM Press, Varanasi, India, 2018, pp. 1–7. `doi:10.1145/3154273.3154315`.
URL `http://dl.acm.org/citation.cfm?doid=3154273.3154315`

[12] M. W. Mahoney, P. Drineas, CUR matrix decompositions for improved data analysis, Proceedings of the National Academy of Sciences 106 (3) (2009) 697–702.

[13] S. Qazi, S. M. Atif, M. B. Kadri, A Novel Compressed Sensing Technique for Traffic Matrix Estimation of Software Defined Cloud Networks, KSII Transactions on Internet and Information Systems 12 (10) (Oct. 2018). `doi:10.3837/tiis.2018.10.004`.
URL `http://itiis.org/digital-library/manuscript/2139`

[14] D. Jiang, X. Wang, L. Guo, H. Ni, Z. Chen, Accurate estimation of large-scale IP traffic matrix, AEU - International Journal of Electronics and Communications 65 (1) (2011) 75–86. `doi:10.1016/j.aeue.2010.02.008`.
URL `https://linkinghub.elsevier.com/retrieve/pii/S1434841110000531`

[15] H. Zhou, L. Tan, Q. Zeng, C. Wu, Traffic matrix estimation: A neural network approach with extended input and expectation maximization iteration, Journal of Network and Computer Applications 60 (2016) 220–232. `doi:10.1016/j.jnca.2015.11.013`.
URL `https://linkinghub.elsevier.com/retrieve/pii/S1084804515002854`

[16] L. Nie, D. Jiang, L. Guo, S. Yu, Traffic matrix prediction and estimation based on deep learning in large-scale IP backbone networks, Journal of Network and Computer Applications 76 (2016) 16–22. `doi:10.1016/j.jnca.2016.10.006`.
URL `https://linkinghub.elsevier.com/retrieve/pii/S1084804516302351`

[17] J. Zhao, H. Qu, J. Zhao, D. Jiang, Towards traffic matrix prediction with LSTM recurrent neural networks, Electronics Letters 54 (9) (2018) 566–568. `doi:10.1049/el.2018.0336`.
URL `https://digital-library.theiet.org/content/journals/10.1049/el.2018.0336`

[18] K. Zhang, Z. Hu, X.-T. Gan, J.-B. Fang, A Network Traffic Prediction Model Based on Quantum-Behaved Particle Swarm Optimization Algorithm and Fuzzy Wavelet Neural Network, Discrete Dynamics in Nature and Society 2016 (2016) 1–11. `doi:10.1155/2016/4135056`.
URL `http://www.hindawi.com/journals/ddns/2016/4135056/`

[19] L. F. Zhang, X. P. Zhang, Network traffic prediction based on BP neural networks optimized by quantum genetic algorithm, Computer Engineering and Science 38 (1) (2016) 114–119.

[20] Y. Lu, H. Li, B. Lu, Y. Zhao, D. Wang, X. Gong, X. Wei, Network Traffic Model with Multi-fractal Discrete Wavelet Transform in Power Telecommunication Access Networks, in: H. Song, D. Jiang (Eds.), Simulation Tools and Techniques, Vol. 295, Springer International Publishing, Cham, 2019, pp. 53–62. `doi:10.1007/978-3-030-32216-8_5`.
URL `http://link.springer.com/10.1007/978-3-030-32216-8_5`

[21] A. Kumar, S. Vidyapu, V. V. Saradhi, V. Tamarapalli, A multi-view subspace learning approach to internet traffic matrix estimation, IEEE Transactions on Network and Service Management (2020).

[22] H. Cai, V. W. Zheng, K. C.-C. Chang, A comprehensive survey of graph embedding: Problems, techniques, and

applications, IEEE Transactions on Knowledge and Data Engineering 30 (9) (2018) 1616–1637.

[23] M. Emami, R. Akbari, R. Javidan, A. Zamani, A new approach for traffic matrix estimation in high load computer networks based on graph embedding and convolutional neural network, Transactions on Emerging Telecommunications Technologies 30 (6) (2019) e3604.

[24] D. Cai, X. He, J. Han, T. S. Huang, Graph regularized nonnegative matrix factorization for data representation, IEEE Transactions on Pattern Analysis and Machine Intelligence 33 (8) (2011) 1548–1560. `doi:10.1109/TPAMI.2010.231`.
URL `http://ieeexplore.ieee.org/document/5674058/`

[25] M. Roughan, Y. Zhang, W. Willinger, L. Qiu, Spatio-Temporal Compressive Sensing and Internet Traffic Matrices (Extended Version), IEEE/ACM Transactions on Networking 20 (3) (2012) 662–676. `doi:10.1109/TNET.2011.2169424`.
URL `http://ieeexplore.ieee.org/document/6058636/`

[26] H.-F. Yu, N. Rao, I. S. Dhillon, Temporal regularized matrix factorization for high-dimensional time series prediction, in: Advances in Neural Information Processing Systems, 2016, pp. 847–855.

[27] V. Leplat, N. Gillis, A. M. Ang, Blind Audio Source Separation With Minimum-Volume Beta-Divergence NMF, IEEE Transactions on Signal Processing 68 (2020) 3400–3410. `doi:10.1109/TSP.2020.2991801`.
URL `https://ieeexplore.ieee.org/document/9084229/`

[28] A. M. S. Ang, N. Gillis, Algorithms and Comparisons of Nonnegative Matrix Factorizations With Volume Regularization for Hyperspectral Unmixing, IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing 12 (12) (2019) 4843–4853. `doi:10.1109/JSTARS.2019.2925098`.
URL `https://ieeexplore.ieee.org/document/8758144/`

[29] R. Chen, H. Li, Online algorithm for foreground detection based on incremental nonnegative matrix factorization, in: 2016 2nd International Conference on Control, Automation and Robotics (ICCAR), Hong Kong, Hong Kong, 2016, pp. 312–317. `doi:10.1109/ICCAR.2016.7486747`.
URL `http://ieeexplore.ieee.org/document/7486747/`

[30] S. M. Atif, S. Qazi, N. Gillis, Improved SVD-based initialization for nonnegative matrix factorization using low-rank correction, Pattern Recognition Letters 122 (2019) 53–59. `doi:10.1016/j.patrec.2019.02.018`.
URL `https://linkinghub.elsevier.com/retrieve/pii/S0167865519300583`

[31] Yin Zhang, Abilene Network Datasets, (Accessed At: 12/05/2020, 03:16:14 PM) (2004).
URL `https://www.cs.utexas.edu/~yzhang/research/AbileneTM/`

[32] S. Uhlig, B. Quoitin, J. Lepropre, S. Balon, Providing public intradomain traffic matrices to the research community, ACM SIGCOMM Computer Communication Review 36 (1) (2006) 83–86.

[33] F. Pompili, N. Gillis, P.-A. Absil, F. Glineur, Two algorithms for orthogonal nonnegative matrix factorization with application to clustering, Neurocomputing 141 (2014) 15–25.

[34] M. Ahookhosh, L. T. K. Hien, N. Gillis, P. Patrinos, Multi-block bregman proximal alternating linearized minimization and its application to sparse orthogonal nonnegative matrix factorization, arXiv preprint arXiv:1908.01402 (2019).

[35] N. Gillis, Nonnegative Matrix Factorization, SIAM, Philadelphia, 2020.
URL `https://my.siam.org/Store/Product/viewproduct/?ProductId=32898069`

[36] Y. Nesterov, A method of solving a convex programming problem with convergence rate $O(1/k^2)$, Soviet Mathematics Doklady 27 (2) (1983) 372–376.

[37] N. Guan, D. Tao, Z. Luo, B. Yuan, NeNMF: An Optimal Gradient Method for Nonnegative Matrix Factorization, IEEE Trans. Signal Process. 60 (6) (2012) 2882–2898. `doi:10.1109/TSP.2012.2190406`.
URL `http://ieeexplore.ieee.org/document/6166359/`

[38] B. O'donoghue, E. Candes, Adaptive restart for accelerated gradient schemes, Foundations of Computational Mathematics 15 (3) (2015) 715–732.

[39] Y. Xu, W. Yin, A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion, SIAM Journal on Imaging Sciences 6 (3) (2013) 1758–1789.

[40] D. Berrar, Cross-validation, Encyclopedia of Bioinformatics and Computational Biology 1 (2018) 542–545.

[41] C. Dorffer, M. Puigt, G. Delmaire, G. Roussel, Fast nonnegative matrix factorization and completion using nesterov iterations, in: International Conference on Latent Variable Analysis and Signal Separation, Springer, 2017, pp. 26–35.

[42] S. Zhang, W. Wang, J. Ford, F. Makedon, Learning from incomplete ratings using non-negative matrix factorization, in: Proceedings of the 2006 SIAM international conference on data mining, SIAM, 2006, pp. 549–553.

[43] Y. Vardi, Network tomography: Estimating source-destination traffic intensities from link data, Journal of the American statistical association 91 (433) (1996) 365–377.

[44] Y. Vardi, Applications of the em algorithm to linear inverse problems with positivity constraints, in: Image models (and their speech model cousins), Springer, 1996, pp. 183–198.

[45] S. Huang, P. Zhao, Y. Ren, T. Li, Z. Xu, Self-paced and soft-weighted nonnegative matrix factorization for data

representation, Knowledge-Based Systems 164 (2019) 29–37.

[46] S. Huang, Z. Xu, Z. Kang, Y. Ren, Regularized nonnegative matrix factorization with adaptive local structure learning, Neurocomputing 382 (2020) 196–209.

[47] P. De Handschutter, N. Gillis, Deep orthogonal matrix factorization as a hierarchical clustering technique, in: 2021 28th European Signal Processing Conference (EUSIPCO), IEEE, 2021.

[48] S. Xu, M. Kodialam, T. Lakshman, S. S. Panwar, Learning based methods for traffic matrix estimation from link measurements, IEEE Open Journal of the Communications Society 2 (2021) 488–499.

[49] L. Nie, Y. Li, X. Kong, Spatio-temporal network traffic estimation and anomaly detection based on convolutional neural network in vehicular ad-hoc networks, IEEE Access 6 (2018) 40168–40176.

[50] S. Huang, Z. Kang, Z. Xu, Auto-weighted multi-view clustering via deep matrix decomposition, Pattern Recognition 97 (2020) 107015.

## CRediT author statement

**Syed Muhammad Atif:** Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Visualization, Investigation, Writing-Original draft preparation. **Nicolas Gillis:** Supervision, Visualization, Writing-Original draft preparation, Writing- Reviewing and Editing. **Sameer Qazi:** Supervision, Visualization, Writing- Reviewing and Editing. **Imran Naseem:** Supervision, Visualization, Writing-Reviewing and Editing.

**Syed Muhammad Atif** completed his Bachelors in Computer Engineering from Sir Syed University of Engineering and Technology. He completed his MS in Computer Networks from Usman Institute of Technology. He is currently pursuing a PhD in Computer Science from PAF-KIET under HEC Indigenous Fellowship Program. He is the author or co-author of 12 publications those were published in reputable journals. His current areas of research are Non Negative Matrix Factorization, Deep Learning, Sentiment Analysis, Network Tomography and Green Routing.

**Nicolas Gillis** received his M.Sc. and Ph.D. degrees in applied mathematics from UCLouvain, Belgium, in 2007 and 2011, respectively. He is currently an associate professor in the Department of Mathematics and Operational Research, University of Mons, Belgium. His research interests include optimization, numerical linear algebra, signal processing, machine learning, and data mining. He received the Householder Award in 2014 and a European Research Council starting grant in 2015. He has published a book on 'Nonnegative Matrix Factorization' in 2020 with SIAM. He currently serves as an associate editor of IEEE Transactions on Signal Processing and SIAM Journal on Matrix Analysis and Applications.

**Sameer Qazi** received his B.E. degree from National University of Sciences and Technology, Pakistan, in 2001 and the MS and PhD degrees from the University of New South Wales, Australia, in 2004 and 2009. He is currently working as Associate Professor in the College of Engineering at Karachi Institute of Economics and Technology, Pakistan. His research interests are Computer Network Optimization Problems, Network Tomography, UAV based video surveillance applications and Cloud Computing.

**Imran Naseem** received the BE degree in electrical engineering from the NED University of Engineering and Technology, Pakistan, in 2002, the MS degree in electrical engineering from the King Fahd University of Petroleum and Minerals (KFUPM), KSA, in 2005, and PhD degree from the University of Western Australia, in 2010. He did his post doctorate in the Institute for Multi-sensor Processing and Content Analysis, Curtin University of Technology, Australia. He joined the College of Engineering, KIET, Pakistan, in 2011 where he is currently an associate professor. He is also an adjunct research fellow in the School of Electrical, Electronic, and Computer Engineering, University of Western Australia. His research interests include pattern classification and machine learning with a special emphasis on biometrics and bioinformatics applications. He has authored several publications in top journals and conferences including the IEEE Transactions on Pattern Analysis and Machine Intelligence, IEEE International Conference on Image Processing etc. His benchmark work on face recognition has received more than 180 citations in less than four years. He is also a reviewer of the IEEE Transactions on Pattern Analysis and Machine Intelligence, the IEEE Transactions on Image Processing, and the IEEE Signal Processing Letters.