

Estimating Network Flow Length Distributions via Bayesian Nonnegative Tensor Factorization

Barış Kurt*

Ali Taylan Cemgil*

Güneş Karabulut Kurt†

Engin Zeydan‡

September 5, 2019

Abstract

In this paper, we develop a framework to estimate network flow length distributions in terms of the number of packets. We model the network flow length data as a three-way array with day-of-week, hour-of-day, flow-length as entities where we observe a count. In a high-speed network, only a sampled version of such an array can be observed and reconstructing the true flow statistics from fewer observations becomes a computational problem.

We formulate the sampling process as matrix multiplication so that any sampling method can be used in our framework as long as its sampling probabilities are written in matrix form. We demonstrate our framework on a high volume real-world data set collected from a mobile network provider with a random packet sampling and a flow-based packet sampling methods. We show that modeling the network data as a tensor improves estimations of the true flow length histogram in both sampling methods.

Keywords: network flow length distribution, packet sampling, flow sampling, nonnegative tensor factorization

*Department of Computer Engineering, Bogazici University, Istanbul, Turkey

†Istanbul Technical University, Istanbul, Turkey

‡Centre Technologic de Telecomunicacions de Catalunya, Castelldefels, Spain, 08860.

Email: engin.zeydan@cttc.cat

1 Introduction

Monitoring network statistics is crucial for the maintenance and infrastructure planning for network service providers. Statistical information about traffic patterns helps a service provider to characterize its network resource usage and user behavior, to infer future traffic demands, to detect traffic and usage anomalies, and to provide insights to improve the performance of the network [1]. However, network monitoring has become a difficult task due to increasingly high-volume and high-speed data over modern networks and in most cases, it requires special hardware. For this reason, sampling [2] becomes a viable approach for extracting statistics from such high-speed networks. In this work, we are interested in one of the most important network statistics, the flow length distribution (FLD).

A network flow is defined as a set of internet protocol (IP) packets with the same signature observed within a limited time period. The flow signature is composed of the IP and port pairs of both source and destination nodes together with level-3 protocol types such as transport control protocol (TCP) or user datagram protocol (UDP). A flow starts with the arrival of the first packet and terminated when the inter-packet timeout is exceeded. The total number of packets in a flow is referred to as the flow length and the length distribution of a set of flows that are terminated in a time window is called flow length distribution.

In this work, we are using one of the most popular methods for collecting per-flow information, i.e., passive measurement. In this method, network packets are processed as they pass through a passive measurement beacon connected to the network, e.g., router. The beacon keeps a look-up table for flow identification. The beacon processes a packet by searching its corresponding flow inside the look-up table using its signature. If such a flow is found, its statistics are updated. Otherwise, the packet is treated as the first packet of a new flow, and the new flow is inserted into the table. Once a flow is terminated, its statistics are transferred to a storage.

The flow length histogram can be calculated exactly by processing every packet that passes through the measurement beacon. In order to implement such a direct method, the monitoring beacon needs to maintain a table to hold information for all active flows on the network. However, substantial amount of concurrent flows with very short packet inter-arrival times of current high-speed networks (on the order of 10 Gbps to 100 Gbps inside carrier's network today) make this brute-force counting method very costly to implement. First of all, this method would require a large amount of memory to record the flow table. Secondly, in a high-speed link, the inter-arrival times between packets, which may be as small as 8 nanoseconds in an OC-768 link,

41 may be smaller than the time required to process flow hash operations such
42 as identifying the packet and updating the flow statistics.

43 The characteristics of the network traffic data inevitably lead to the de-
44 velopment of alternative methods for measurement such as random sampling,
45 where a fraction of the network traffic is randomly selected and processed.
46 The simplest sampling method is the uniform packet sampling [3, 4, 5, 6],
47 used in commercial systems [7] and [8]. In uniform sampling, each packet
48 is selected with a predefined constant probability. This approach is easy
49 to implement since it does not require the flow identification of each packet.
50 However, recovering the true flow length distribution from the random packet
51 sampled traffic is a challenging problem. The unbiased estimator of the orig-
52 inal flow length n for sampling probability p is $\hat{n}(m) = m/p$, where m
53 is observed flow length. The relative error of this estimator, calculated as
54 $\sqrt{1-(p/n)^2}$ [3], grows unboundedly for short flows as the sampling rate
55 gets smaller. The high error on the small flow lengths comes from the fact
56 that most of the samples are collected from longer flows.

57 Flow-based adaptive sampling methods [9, 10, 11, 12, 13, 14] were pro-
58 posed for more accurate flow length estimation. In these methods, each
59 incoming packet is processed and then sampled with a probability that is a
60 function of the current sampled length of the flow that the packet belongs
61 to. Here, the main idea is to use a smaller memory by compressing the
62 flow statistics counters in the router. However, these methods need to be
63 implemented on specialized -and expensive- hardware due to the mandatory
64 packet identification and lookup step.

65 Both packet-based and flow-based adaptive sampling methods rely on
66 numerical methods to recover the true FLD. In this work, we propose a
67 framework that can be used to recover the true FLD from the sampled ob-
68 servation obtained by any sampling method. This framework uses a variant
69 of the nonnegative tensor factorization (NTF) model, which we call the thin
70 nonnegative tensor factorization (ThinNTF), where the "thin" prefix em-
71 phasizes that the factorization is applied directly to the sampled, or namely
72 "thinned" data.

73 In our framework, the network traffic data is modeled as a 3-way array,
74 containing the number of flow length observations, with dimensions inter-
75 preted as 1) flow length, 2) hour-of-day and 3) day-of-week to capture the
76 hourly and daily periodicity in the data. The nonnegative factorization of
77 this tensor basically gives us estimates in the form of a non-parametric mix-
78 ture model. Therefore our model is an improvement of the non-parametric
79 flow length models in [3] and [6] by having the capability of modeling data
80 with an arbitrary amount of mixture components and use the periodicity.

81 While the ordinary NTF model [15] factorizes an observation tensor, the

82 ThinNTF directly factorizes its sampled version and recovers the original
83 tensor from the estimated factors. We take a fully Bayesian approach here
84 and provide a generative model for the ThinNTF and a variational Bayes
85 algorithm for inference. The contributions in this paper can be listed as:

- 86 • We model one week of flow length observations as a 3-dimensional
87 tensor and observe the periodic behavior.
- 88 • We propose a novel tensor factorization scheme, ThinNTF, which is
89 able to find the factors of a latent tensor from its sampled counterpart.
90 By doing so, we also solve the reconstruction problem.
- 91 • We apply ThinNTF to real-world data sampled with two different sam-
92 pling methods: uniform random packet sampling and flow-based adap-
93 tive sampling.

94 The structure of the paper is as follows. In Section 2, we provide the
95 related works on network sampling and tensor factorization. In Section 3, we
96 describe our real-world data and how we visualize it as a tensor. Addition-
97 ally, we describe the sampling methods that we used to sample the data. In
98 Section 4, we describe our ThinNTF model and the variational Bayes algo-
99 rithm for estimating the factors. In Section 5, we describe our real-world data
100 collection architecture. In section 6, we present our synthetic and real-world
101 experiments and results. Finally, in Section 7, we draw our conclusions.

102 2 Related Work

103 Sampling methods have long been applied to network traffic monitoring. A
104 survey on fundamental network sampling strategies is given in [2]. Uniform
105 packet sampling is extensively studied by various authors. Duffield et al. [3]
106 propose the first non-parametric model for flow length distribution and pro-
107 vides a maximum likelihood estimation to recover the flow lengths. Riberio et
108 al. [4] show that using protocol specific information gives better flow length
109 distribution estimates in TCP flows. Yang et al. [6] adopt the maximum
110 likelihood approach to estimate both flow length and flow volume (number
111 of bytes in a flow) distributions. Additionally, they model the data with a
112 non-parametric mixture model of two components, where the first component
113 models small flows and the second models large ones.

114 Flow-based sampling methods are proposed as alternatives to the uni-
115 form packet sampling since packet sampling has theoretical limitations when
116 recovering true flow statistics [5]. These methods process every incoming

117 packet and apply sampling conditionally. Kumar et al. [16, 13] propose two
 118 different algorithms where the flow size counters are compressed statistically.
 119 They also propose a non-uniform packet sampling algorithm based on sketch
 120 counting [12]. Hu et al. [10, 14] propose another non-uniform packet sam-
 121 pling algorithm, called adaptive nonlinear sampling (ANLS) for estimating
 122 flow lengths per each flow, and then adopts this method to flow volume [11].
 123 In our experiments, we are going to use ANLS as an example of flow-based
 124 sampling methods since it is the current state-of-the-art non-uniform sam-
 125 pling method.

126 Nonnegative Tensor Factorization is the generalization of the nonnegative
 127 Matrix factorization (NMF) [17] to multiway arrays. In NMF, a nonnegative
 128 matrix is approximated with a multiplication of two nonnegative matrices.
 129 Minimizing the Kullback-Leibler divergence between the initial matrix and
 130 multiplied factors is a popular formulation of this method and can be solved
 131 with fixed-point iterations [18] or full Bayesian methods [19]. NMF has been
 132 used in many applications such as spectral data analysis [20], face recognition
 133 [17] and document clustering [21].

134 Modeling the flow length distribution as a mixture of distributions is
 135 first proposed by [6]. However, according to our best knowledge, there is no
 136 previous work that models a large volume of flow length data as a tensor.
 137 This work fills a gap in the literature by introducing tensor factorization
 138 methodology to network monitoring.

139 **3 Problem Description**

140 We describe our problem as a tensor thinning problem, where the count
 141 entities of the original flow lengths are stored in a tensor. We formulate the
 142 sampling process as a matrix multiplication operated on this data tensor. In
 143 order to do that, each sampling model should be represented as a matrix that
 144 transforms the original data tensor to a sampled one. We provide matrices
 145 for two sampling models: uniform packet sampling and ANLS flow-based
 146 packet sampling.

147 **3.1 Notation and Indexes**

148 For a clear notation, the scalar values are denoted by lightface letters, such as
 149 the index variable j and its maximum value J . The vectors are represented
 150 by boldface lower case letters, such as vector \mathbf{x} . Boldface upper case letters
 151 represent matrices, such as \mathbf{F} ; \mathbf{H} and \mathbf{D} , and the tensors are represented with
 152 calligraphic upper case letters i.e \mathcal{X} . The individual entries in matrices and

Index	Range	Description
i	$[1; I]$	Original flow lengths
	$[0; I]$	Sampled flow lengths
j	$[1; J]$	Hours of day
k	$[1; K]$	Days
r	$[1; R]$	Components

Table 1: Indexes in the model.

153 tensors are written like scalars, i.e. $f_{i,r}$ and $x_{i,j,k}$. The index $:$ denotes all the
154 entries in the given dimension. For example $S_{i,:}$ is the i^{th} row of the \mathbf{S} matrix
155 and $\mathcal{X}_{i,:,:}$ is the i^{th} slice of the tensor \mathcal{X} in the first dimension.

156 The index parameters are also fixed for clarity. The list of indexes, their
157 ranges and semantic descriptions are given in Table 1. For example, the i
158 index always represents an original flow length, while i presents a sampled
159 flow length. The range of i starts from 0, since all of the packets of a flow
160 may be discarded during the sampling process yielding a zero-length sampled
161 flow, which is never observed.

162 3.2 Data Tensor

163 The original flow length data is represented in an $I \times J \times K$ tensor \mathcal{X} , with
164 individual elements $x_{i,j,k}$, regarded as the number of flows that has length
165 i measured at the hour j of the day k . In this setup, I is the maximum
166 flow length value, J is the hours of the day and K is the days of the week.
167 For our real-world data, collected continuously for 1 week, these values are
168 $I = 2000000; J = 24$ and $K = 7$.

169 Working with large maximum flow size is not feasible for two reasons.
170 First, learning a mixture model where each flow component has 2 million
171 parameters is not a good formulation for this problem. Secondly, 99.9% of
172 flows in our data have less than 100 packets, which means the tensor \mathcal{X} will
173 be very sparse for $i > 100$. The clamping process can be defined as

$$\bar{\mathcal{X}}_{i,j,k} = \begin{cases} \mathcal{X}_{i,j,k} & \text{for } i < I_{max} \\ \sum_{l=I_{max}}^I \mathcal{X}_{l,j,k} & \text{for } i = I_{max} \end{cases} \quad (1)$$

174 where $\bar{\mathcal{X}}$ is the clamped tensor. The clamping does not require any change
175 in the model and inference equations that are given in Section 4. Therefore,
176 for notational clarity we only use \mathcal{X} as the generic original data tensor.

177 Figure 1 shows the vertical slices of our unsampled real-world data tensor
178 \mathcal{X} , collected at the backbone of a mobile operator during a one week period,

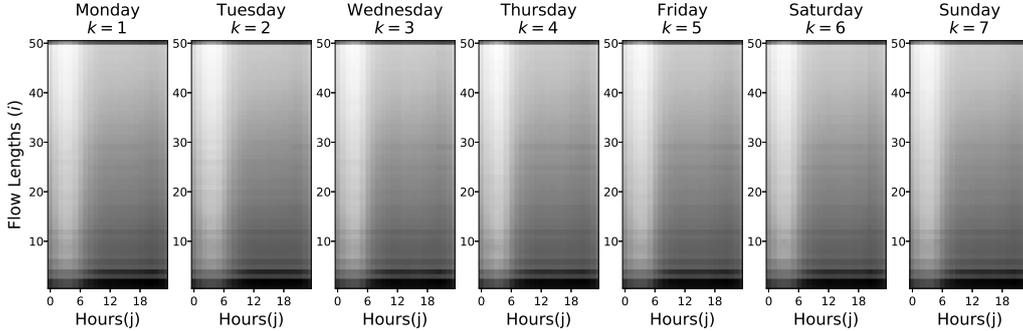


Figure 1: Slices of the original flow length tensor \mathcal{X} .

179 from Monday to Sunday, and clamped at $l_{max} = 50$. The intensity images
 180 are generated from the logarithm of the flow length counts. The daily and
 181 hourly patterns are easily recognizable in the original FLD data.

182 3.3 Sampling Methods

183 Independent of the sampling method, we can define an $l \times (l + 1)$ size \mathbf{S}
 184 matrix, where l is the maximum flow length with entries $S_{l,i}$ interpreted
 185 as the probability of sampling packets from an original flow of length l .
 186 Naturally, \mathbf{S} is a lower diagonal matrix of the form

$$\mathbf{S} = \begin{bmatrix} S_{1,0} & S_{1,1} & 0 & 0 & \dots & 0 \\ S_{2,0} & S_{2,1} & S_{2,2} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ S_{l,0} & S_{l,1} & S_{l,2} & S_{l,3} & \dots & S_{l,l} \end{bmatrix} \quad (2)$$

187 where its l^{th} row defines a probability distribution for the sampled flow
 188 length of a flow of size l . Given a flow size distribution $\mathbf{x} \in \mathbb{Z}^l$, where \mathbb{Z} is
 189 the set of nonnegative integers, the expected sampled flow length distribution
 190 would be given by $\hat{\mathbf{y}} = \mathbf{S}^T \mathbf{x}$. It immediately follows that the sampled flow size
 191 distribution y has length $l + 1$, with y_0 is the proportion of sampled flows with
 192 none of their packets sampled. During the sampling process, this value will
 193 never be observed naturally since the flow identification is performed only on
 194 selected packets. In all experiments throughout this paper, the \mathbf{y} vector (or
 195 \mathcal{Y} tensor, which will be described later on) will be element-wise multiplied
 196 with a binary mask vector \mathbf{m} (or binary mask tensor \mathcal{M}), whose entries are
 197 set to 1 except the ones corresponding to zero sampled flow lengths, in order
 198 to simulate the real life scenario.

Algorithm 1 Uniform Packet Sampling Algorithm

```
1: function SAMPLEUNIFORMLYRANDOM( , flow_table, packet)
2:   if > rand_double(0, 1) then
3:     flow = flow_table.lookup(packet)
4:     if flow is null then
5:       flow = new Flow(packet)
6:     else
7:       flow.length += 1
8:     flow_table.insert_or_update(flow)
```

199 For any given sampling method, we can calculate the \mathbf{S} matrix directly
200 if a closed-form expression is available. Otherwise, it can be approximated
201 by simulating the sampling process and counting the sampling statistics. In
202 this paper, both uniform random sampling and the ANLS provide closed-
203 form expressions for the calculation of \mathbf{S} matrix.

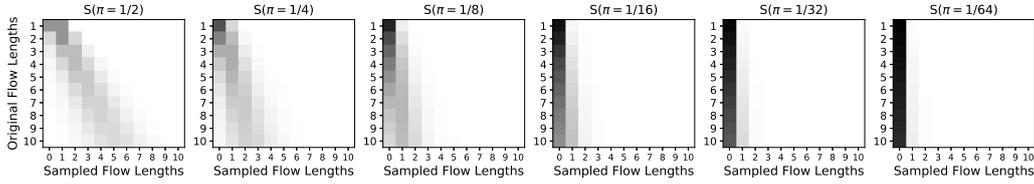
204 An important practical issue is that, if the original tensor \mathcal{X} is clamped
205 at l_{max} , the \mathbf{S} matrix must also be clamped. In that case a last row entry
206 $S_{l_{max},i}$ must present the probability of selecting i packets from a flow of length
207 greater or equal to l_{max} . This clamping operation can be done by calculating
208 a full size \mathbf{S} matrix first, and setting $S_{l_{max},i} \propto \sum_{j=l_{max}}^i S_{j,i}$ with a naive
209 assumption that after l_{max} the original flow sizes are uniformly distributed.

210 3.3.1 Uniform Sampling Method

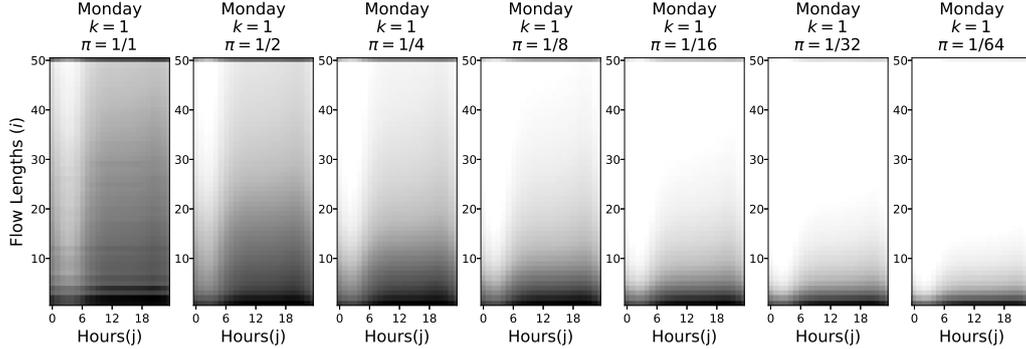
211 In uniform sampling, each packet is processed with a fixed probability of p ,
212 irrespective of the flow it belongs to. In this method, the sampling matrix
213 entries $S_{j,i}$ are calculated directly through Binomial distribution with i trials
214 and p success probability, i.e.,

$$S_{j,i} = \begin{cases} \binom{i}{j} (1-p)^{i-j} p^j & \text{for } j \leq i \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

215 Algorithm 1 describes how the flow table is updated with uniform sam-
216 pling upon the arrival of a new packet. The algorithm uniformly draws a
217 random number in interval $[0;1]$ and if it is less than p , processes the packet,
218 otherwise the packet is discarded. For processing the packet, a look-up op-
219 eration is performed on the flow table to find and update the flow that the
220 packet belongs to. If no such flow is found, a new flow is created using the
221 packet's signature.



(a) Uniform sampling matrices.



(b) Uniformly sampled Monday data with varying sampling probabilities.

Figure 2: Sampling matrices for two different sampling schemes.

222 Figure 2a shows the top 10×11 entries of the lower diagonal \mathbf{S} matrices
 223 with different sampling probabilities. As the sampling probability gets
 224 smaller, fewer packets from a flow gets observed, and the flow may even
 225 be missed when none of its packets are observed. The rightmost sampling
 226 matrix shows the case when $\pi = 1/64$, where the matrix has a very high
 227 concentration of zero-length sampled flows.

228 Figure 2b shows the original Monday data (the leftmost matrix) and its
 229 sampled versions under uniform sampling with the probabilities shown on
 230 the top sampling matrices. Here we see that for $\pi = 1/64$, the observed flow
 231 lengths are mostly less than 10, while the majority are not observed at all.

232 3.3.2 ANLS Sampling Method

233 The ANLS [10] will be used as the representative of the flow-based adaptive
 234 sampling methods. In ANLS, each packet is sampled according to the number
 235 of packets previously sampled from its corresponding flow. If a sampled flow
 236 has length x , the probability of its next packet to be sampled ($\rho(x; u)$) is
 237 calculated as

Algorithm 2 ANLS Sampling Algorithm

```
1: function SETUPANLS( $u$ )
2:    $f[0] = 0$ 
3:   for  $i \in [1; I]$  do
4:      $f[i] = ((1 + u)^i - 1) = u$ 
5:      $\rho[i] = 1 = (f[i - 1] - f[i])$ 
6:   return  $f, p$ 
7: function SAMPLEWITHANLS( $p, \text{flow\_table}, \text{packet}$ )
8:    $\text{flow} = \text{flow\_table.lookup}(\text{packet})$ 
9:   if  $\text{flow}$  is null then
10:     $\text{flow} = \text{new Flow}(\text{packet})$ 
11:   else if  $p[\text{flow.length}] > \text{rand\_double}(0, 1)$  then
12:     $\text{flow.length} += 1$ 
13:    $\text{flow\_table.insert\_or\_update}(\text{flow})$ 
```

$$f(x; u) = [(1 + u)^x - 1] = u \quad (4)$$

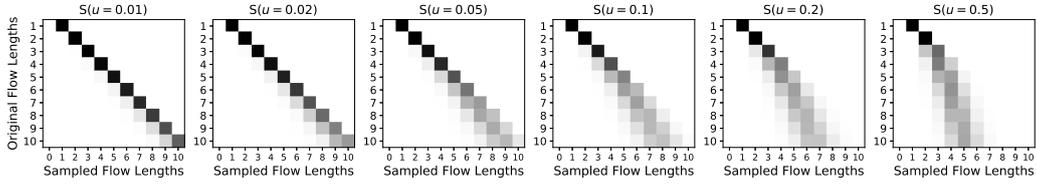
$$\rho(x; u) = 1 = [f(x - 1; u) - f(x; u)] \quad (5)$$

238 Here, $f(x; u)$ is a monotonically increasing function of flow length x ,
239 parametrized with u , which makes $\rho(x; u)$ monotonically decreasing. The u
240 parameter determines the tendency of sampling packets. As u gets smaller,
241 more packets are sampled and estimating original flow lengths gets easier.

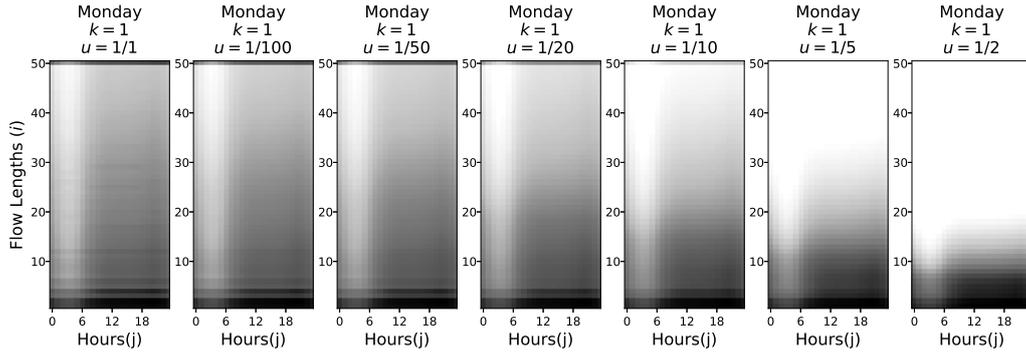
242 The ANLS method is described in details in Algorithm 2. Prior to the
243 sampling, the \mathbf{f} and \mathbf{p} vectors are calculated in the SetupANLS function,
244 according to equations 4 and 5. During sampling, for each incoming packet
245 a look-up operation is performed unconditionally. If the corresponding flow
246 is found, it is updated with probability relative to its current observed size.
247 Otherwise, a new flow is created, ensuring that every flow is observed with
248 at least one packet.

249 We calculate the sampling matrix \mathbf{S} recursively for ANLS. In this method,
250 the first packet is always sampled since $\rho(1; u) = 1$ independent of u . We
251 start by assigning all zero sampling probabilities as $s_{:,0} = 0$ and $s_{1,1} = 1$.
252 The recursive equation for calculating the sampling matrix can be deduced
253 as

$$s_{i,:} \propto s_{i-1,:} \cdot \rho(i-1; u) + s_{i-1,:} \cdot (1 - \rho(i-1; u)) \quad (6)$$



(a) ANLS sampling matrices.



(b) ANLS sampling with varying u .

Figure 3: The visualization of the Monday slice with uniform and ANLS sampling methods.

254 Figure 3 shows the ANLS sampling matrices for first 10 flow lengths
 255 and the Monday data sampled with them respectively, similarly to Figure 2.
 256 First, we can see that when u is small, the \mathbf{S} matrix looks like identity and
 257 as u gets larger, the sampling probability of large flows decreases. Secondly,
 258 compared to uniform sampling, the ANLS method has much higher sampling
 259 ratios than uniform sampling. However, operating with such high sampling
 260 ratios would require specialized hardware in real time.

261 4 Methodology

262 Our methodology is based on the nonnegative factorization of the data tensor.
 263 Our model, which we call ThinNTF, introduces the sampling matrix as a
 264 constant factor to the original NTF with the Poisson-Gamma observation
 265 model. The rationale for using factorization for recovering true flow sizes
 266 is that the flow size distributions have daily periodic behavior, as we show
 267 in Section 3. Inferring the factors, instead of the original matrix, requires
 268 less parameter estimation and results in smoother estimates compared to the
 269 standard maximum likelihood estimation described in [3]. In this section, we
 270 first describe the original tensor factorization model and we provide our novel

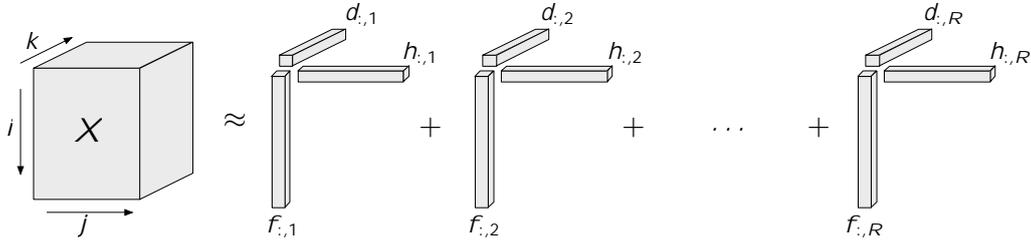


Figure 4: PARAFAC Factorization

271 version: the ThinNTF. At the end of the section, we present the full-Bayesian
 272 variational Bayes algorithm for the inference.

273 4.1 Nonnegative Tensor Factorization

274 NTF is the generalization of the 2-dimensional NMF model to multiple di-
 275 mensions. In NTF, an N-dimensional tensor is approximated by the multi-
 276 plication of lower dimensional factors. Unlike NMF, tensor factorization can
 277 be done in multiple ways. In this work, we are going to use the PARAFAC
 278 [22, 23, 24] factorization scheme. In PARAFAC, an $I_1 \times I_2 \times \dots \times I_N$ tensor
 279 is approximated by $I_n \times R$ matrices for $n \in [1; N]$. Here, R is the number
 280 of components, i.e. the number of clusters in the data. Figure 4 shows the
 281 PARAFAC factorization of our FLD tensor \mathcal{X} , into 3 factors: an $I \times R$ factor
 282 \mathbf{F} for representing the flow length clusters, a $J \times R$ factor \mathbf{H} for representing
 283 hourly behavior and a $K \times R$ factor \mathbf{D} for representing the daily behavior of
 284 the data. Every single entry of the \mathcal{X} tensor is approximated by

$$x_{i,j,k} \approx \hat{x}_{i,j,k} = \sum_r f_{i,r} h_{j,r} d_{k,r} \quad (7)$$

285 Bro [25] explains that the PARAFAC factorization is unique under cer-
 286 tain circumstances, where uniqueness is defined as being unable to rotate the
 287 factorization without loss of fit. NMF and NTF are statistical models that
 288 impose nonnegativity constraint without uniqueness property. The unique-
 289 ness may be important if individual factors are of special interest. In our
 290 case, we are concerned with the estimation of the original data tensor \mathcal{X}
 291 from sampled tensor \mathcal{Y} , but not the individual factors for any interpretation.
 292 Our problem is more close to a missing value imputation problem, hence
 293 uniqueness is not a requirement.

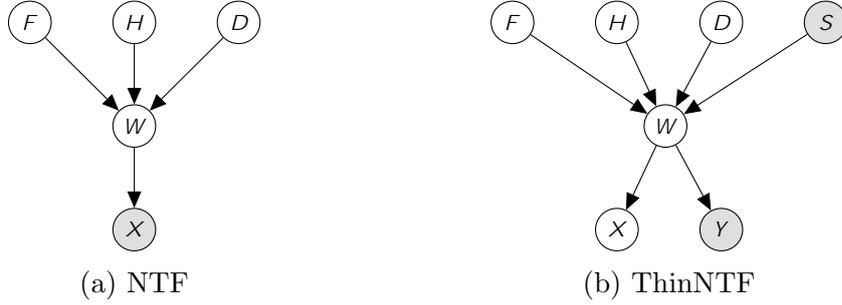


Figure 5: Graphical models representing the dependency structure of NTF and ThinNTF models in PARAFAC scheme.

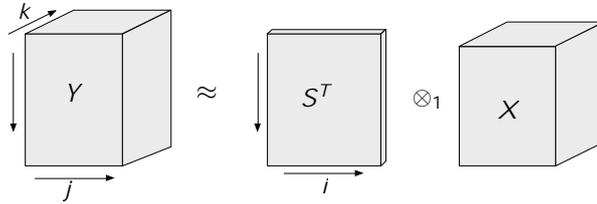


Figure 6: ThinNTF Model

294 4.2 Thin Nonnegative Tensor Factorization

295 ThinNTF is basically an NTF with an additional constant factor, which in
 296 our case is the sampling matrix \mathbf{S} . Figure 5 shows the graphical models
 297 of the NTF and the ThinNTF models for factorizing original and sampled
 298 flow length observations. In the graphical models, the shaded nodes are the
 299 observed entities, and the unshaded ones are the latent entities.

300 In Section 3.3, we have described the sampling process as a matrix mul-
 301 tiplication operation with a sampling matrix \mathbf{S} . In ThinNTF, this sampling
 302 matrix operates on the original tensor \mathcal{X} and creates a thinned version of
 303 it, which we call \mathcal{Y} , by down-sampling its entries according to a sampling
 304 scheme, as shown in Figure 6. The entries of \mathcal{Y} tensor $y_{j:k}$ presents the
 305 number of flows of sampled-length τ , at hour j at day k . The \otimes_1 operation
 306 denotes the 1-mode product of matrix \mathbf{S}^T and tensor \mathcal{X} , which corresponds
 307 to the set of matrix multiplications $\mathcal{Y}_{:::k} = \mathbf{S}^T \mathcal{X}_{:::k}$ for $k \in [1; K]$.

308 In this scheme, one can immediately suspect that \mathcal{X} can be estimated by
 309 $(\mathbf{S}^T)^{-1} \otimes_1 \mathcal{Y}$. However, this solution is not feasible for several reasons. First,
 310 the \mathbf{S} matrix is not square, hence not invertible. Instead its pseudo-inverse
 311 can be calculated but this does not impose nonnegativity. Moreover, the top
 312 slice of the \mathcal{Y} tensor, which stores the number of flows with zero-sampled
 313 size is never observed, hence must be estimated. Therefore we need a solid

Tensor	Index Set	Description
\mathcal{X}	$i;j;k$	Original flow length tensor
\mathcal{Y}	$;j;k$	Sampled flow length tensor
\mathcal{M}	$;j;k$	Mask tensor
\mathcal{W}	$;i;j;k;r$	Latent variable tensor
\mathbf{F}	$i;r$	Flow length factor
\mathbf{H}	$j;r$	Hour of day factor
\mathbf{D}	$k;r$	Day of week factor
\mathbf{S}	$i;$	Sampling matrix
$\mathbf{A}^F; \mathbf{B}^F$	$i;r$	Gamma priors for \mathbf{F}
$\mathbf{A}^H; \mathbf{B}^H$	$j;r$	Gamma priors for \mathbf{H}
$\mathbf{A}^D; \mathbf{B}^D$	$k;r$	Gamma priors for \mathbf{D}

Table 2: Tensors in the model and their corresponding index sets.

314 statistical model and an inference method to estimate \mathcal{X} under this model.

315 In ThinNTF, we observe the \mathcal{Y} tensor, but try to factorize the \mathcal{X} tensor,
316 which is latent (Figure 5b). In the end, the factors of \mathcal{X} are going to provide
317 us an approximation $\hat{\mathcal{X}}$ which solves the original flow length distribution
318 reconstruction problem. We mathematically express this approximation as

$$y_{j;k} \approx \hat{y}_{j;k} = \sum_{i;r} s_{;i} f_{i;r} h_{j;r} d_{k;r} \quad (8)$$

319 where $\mathbf{F}; \mathbf{H}$ and \mathbf{D} are described in exactly the same way in the original NTF
320 case. In subsections 3.3.1 and 3.3.2, we described two different \mathbf{S} matrices for
321 two different schemes. ThinNTF model can be employed with any sampling
322 method as long as it is described with a sampling matrix.

323 4.3 Generative Model

324 Taking Bayesian approach, we first provide a generative model for the Thin-
325 NTF, then describe how we can estimate the posterior probabilities of model
326 parameters (in this case, the factor matrices) conditioned on the sampled
327 flow length observations \mathcal{Y} and the sampling matrix \mathbf{S} using the well known
328 Bayes rule. Table 2 contains all tensors and matrices used in the model
329 together with their index sets.

The original and latent data tensor \mathcal{X} , and the sampled and observed data tensor Y have nonnegative integer entries. The natural probability distribution for this type of count data is the Poisson distribution. We assume

that each entry of a latent 5-dimensional tensor \mathcal{W} is drawn from a Poisson distribution whose parameters are functions of sampling matrix \mathbf{S} and factors \mathbf{F} ; \mathbf{H} and, \mathbf{D} , such as

$$w_{:ij:k:r} \sim \mathcal{PO}(w_{:ij:k:r}; s_{:i}f_{i,r}h_{j,r}d_{k,r}) \quad (9)$$

where the Poisson distribution is defined as

$$\mathcal{PO}(w; \lambda) = \exp(w \log \lambda - \lambda - \log \Gamma(w + 1)) \quad (10)$$

We choose the prior distributions for the factor entries as the Gamma distribution since it is the conjugate prior of Poisson distribution [26]. For each entry of factor \mathbf{F} , we write

$$f_{i,r} \sim \mathcal{G} \left(f_{i,r}; a_{i,r}^f, \frac{b_{i,r}^f}{a_{i,r}^f} \right) \quad (11)$$

330 where the Gamma distribution is described as

$$\mathcal{G}(f; \lambda; \Theta) = \exp((\lambda - 1) \log f - \frac{f}{\Theta} - \log \Theta - \log \Gamma(\lambda)) \quad (12)$$

331 with shape parameter λ and scale parameter Θ . In our generative model,
 332 the parameters for Gamma distributions are $\lambda = a_{i,r}^f$ and $\Theta = b_{i,r}^f = a_{i,r}^f$ respec-
 333 tively. This means that the mean of $f_{i,r}$ is $\Theta = b_{i,r}^f$, which is independent
 334 of $a_{i,r}^f$. The variance of $f_{i,r}$ becomes $\Theta^2 = (b_{i,r}^f)^2 = a_{i,r}^f$, which means that as
 335 $a_{i,r}^f$ gets smaller, the factors gets sparser.

336 In order to avoid repetition, we are going to omit the equations regarding
 337 the factors \mathbf{H} and \mathbf{D} throughout the paper. These factors behave exactly
 338 like factor \mathbf{F} and it's easy to derive equations related to these factors once
 339 their corresponding equation for \mathbf{F} is given.

340 Finally, we generate \mathcal{X} and \mathcal{Y} tensor from \mathcal{W} . Each entry $w_{:ij:k:r}$ of \mathcal{W}
 341 can be interpreted as the number of original flows of length i , generated on
 342 hour j , day k , by cluster r and observed as length \cdot . By summing \mathcal{W} over
 343 dimensions cluster (r) and original lengths (i), we get the sampled observa-
 344 tions tensor \mathcal{Y} . Similarly, by summing \mathcal{W} over dimensions cluster (r) and
 345 sampled lengths (\cdot), we get the original flow length tensor \mathcal{X} . The whole
 346 generative process is summarized in Algorithm 3. The set of all indexes and
 347 tensors in the model are summarized in Tables 1 and 2 respectively.

Algorithm 3 ThinNTF Generative Model

```
1: function RANDINIT( $\mathbf{S}; \mathbf{A}^F; \mathbf{B}^F; \mathbf{A}^H; \mathbf{B}^H; \mathbf{A}^D; \mathbf{B}^D$ )
    // Sample factor  $\mathbf{F}$  from  $\text{Gamma}(\mathbf{A}^F; \mathbf{B}^F)$ 
2:   for all  $i \in [1; I]; r \in [1; R]$  do
3:      $f_{i,r} \sim \mathcal{G}\left(f_{i,r}; a_{i,r}^f; \frac{b_{i,r}^f}{a_{i,r}^f}\right)$ 

    // Sample factor  $\mathbf{H}$  from  $\text{Gamma}(\mathbf{A}^H; \mathbf{B}^H)$ 
4:   for all  $k \in [1; K]; r \in [1; R]$  do
5:      $h_{j,r} \sim \mathcal{G}\left(h_{j,r}; a_{j,r}^h; \frac{b_{j,r}^h}{a_{j,r}^h}\right)$ 

    // Sample factor  $\mathbf{D}$  from  $\text{Gamma}(\mathbf{A}^D; \mathbf{B}^D)$ 
6:   for all  $j \in [1; J]; r \in [1; R]$  do
7:      $d_{k,r} \sim \mathcal{G}\left(d_{k,r}; a_{k,r}^d; \frac{b_{k,r}^d}{a_{k,r}^d}\right)$ 

    // Sample latent tensor  $\mathcal{W}$  from Poisson distributions
8:   for all  $l \in [1; I+1]; i \in [1; I]; j \in [1; J]; k \in [1; K]; r \in [1; R]$  do
9:      $w_{i,j,k;r} \sim \mathcal{PO}(w_{i,j,k;r}; s_{i,j,k;r} f_{i,r} h_{j,r} d_{k,r})$ 
10:  return  $\{\mathbf{F}, \mathbf{H}, \mathbf{D}, \mathcal{W}\}$ 

11: function GENERATEDATA( $\mathbf{S}; \mathbf{A}^F; \mathbf{B}^F; \mathbf{A}^H; \mathbf{B}^H; \mathbf{A}^D; \mathbf{B}^D$ )
    // Randomly initialize factors and latent tensor
12:   $\{\mathbf{F}; \mathbf{H}; \mathbf{D}; \mathcal{W}\} \leftarrow \text{RANDINIT}(\mathbf{S}; \mathbf{A}^F; \mathbf{B}^F; \mathbf{A}^H; \mathbf{B}^H; \mathbf{A}^D; \mathbf{B}^D)$ 
    // Generate original tensor  $\mathcal{X}$ 
13:  for all  $i \in [1; I]; j \in [1; J]; k \in [1; K]$  do
14:     $x_{i,j,k} = \sum_{r} w_{i,j,k;r}$ 

    // Generate sampled tensor  $\mathcal{Y}$ 
15:  for all  $l \in [1; I+1]; j \in [1; J]; k \in [1; K]$  do
16:     $y_{j,k} = \sum_{i,r} w_{i,j,k;r}$ 
17:  return  $\{\mathbf{F}, \mathbf{H}, \mathbf{D}, \mathcal{W}, \mathcal{X}, \mathcal{Y}\}$ 
```

348 4.4 Variational Bayes

349 After defining the generative model, we can infer the factors \mathbf{F} , \mathbf{H} , and \mathbf{D}
350 of a sampled flow length observation tensor \mathcal{Y} . In the original NMF paper,
351 Lee and Seung [17] provide fixed-point update equations for inferring the
352 factors. Bro [25] gives similar fixed-point equations for updating the factors in
353 PARAFAC factorization. Cemgil [19] shows that these updates correspond to

354 the Kullback-Leibler minimization between the original matrix (or tensor \mathcal{X})
 355 and the approximated one ($\hat{\mathcal{X}}$), and also provides a full Bayesian variational
 356 algorithm for the matrix factorization. Ermis et. al. [15] provide a similar
 357 variational algorithm for the Gamma-Poisson tensor factorization.

We start our Bayesian inference by calculating the posterior distributions over the factors $\mathbf{F}; \mathbf{H}$ and \mathbf{D} conditioned on observed tensor \mathcal{Y} . For notational clarity, we introduce $\theta = (\mathbf{A}^F; \mathbf{B}^F; \mathbf{A}^H; \mathbf{B}^H; \mathbf{A}^D; \mathbf{B}^D)$ as the list of model hyper-parameters. The log-likelihood observing \mathcal{Y} under the model parameters θ is written as

$$\log p(\mathcal{Y} | \theta; \mathbf{S}) = \log \int_{\mathbf{F}, \mathbf{H}, \mathbf{D}} d\mathbf{F} d\mathbf{H} d\mathbf{D} \sum_{\mathcal{W}} p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \theta; \mathbf{S}) \quad (13)$$

This log-likelihood is intractable due to the integration over the latent factors, but it is lower bounded as

$$\log p(\mathcal{Y} | \theta; \mathbf{S}) \leq \mathcal{L} \quad (14)$$

$$= \langle \log p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \theta; \mathbf{S}) \rangle_{q(\mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D})} + \mathcal{H}_{q(\mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D})} \quad (15)$$

where q is an auxiliary joint distribution of latent factors. This bound is tight when $q(\mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D}) = p(\mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \mathcal{Y}; \theta; \mathbf{S})$. However, this is also intractable to calculate. Instead, we use a variational approximation [27] for q such that

$$q(\mathcal{W}) \propto \exp \left(\langle \log p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \theta) \rangle_{q(\mathbf{F}, \mathbf{H}, \mathbf{D})} \right) \quad (16)$$

$$q(\mathbf{F}) \propto \exp \left(\langle \log p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \theta) \rangle_{q(\mathcal{W}, \mathbf{H}, \mathbf{D})} \right) \quad (17)$$

$$q(\mathbf{H}) \propto \exp \left(\langle \log p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \theta) \rangle_{q(\mathcal{W}, \mathbf{F}, \mathbf{D})} \right) \quad (18)$$

$$q(\mathbf{D}) \propto \exp \left(\langle \log p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \theta) \rangle_{q(\mathcal{W}, \mathbf{F}, \mathbf{H})} \right) \quad (19)$$

358 where we iteratively update the posterior distribution of each factor by calcu-
 359 lating the expectation of the logarithm of the full joint likelihood $p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D})$
 360 under the posteriors of all other latent factors.

361 4.5 Update Equations

Here we provide the update equations for $q(\mathcal{W})$ and $q(\mathbf{F})$. The updates of $q(\mathbf{H})$ and $q(\mathbf{D})$ can be easily deduced from the update equations of $q(\mathbf{F})$.

The full joint likelihood whose expectation will be calculated at each step is

$$\mathcal{J} = \log \rho(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D} | \cdot) \quad (20)$$

$$\begin{aligned} &= \log \rho(\mathcal{Y} | \mathcal{W}) + \log \rho(\mathcal{W} | \mathbf{F}; \mathbf{H}; \mathbf{D}) \\ &\quad + \log \rho(\mathbf{F} | \cdot) + \log \rho(\mathbf{H} | \cdot) + \log \rho(\mathbf{D} | \cdot) \end{aligned} \quad (21)$$

where $\mathcal{Y} | \mathcal{W}$ is a degenerate distribution (()) to make sure the summation of $\sum_{i;r} \mathcal{W}$ equals \mathcal{Y} . By inserting the necessary Poisson and Gamma distributions given in the generative model into the Equation 21 we get the following expression

$$\begin{aligned} \mathcal{J} &= \sum_{j;k} m_{j;k} \log \left(y_{j;k} - \sum_{i;r} w_{i;j;k;r} \right) \\ &\quad + \sum_{i;r} \sum_{j;k} m_{j;k} \left(w_{i;j;k;r} \log s_{i;j;k;r} - s_{i;j;k;r} - \log \Gamma(w_{i;j;k;r} + 1) \right) \\ &\quad + \sum_{i;r} \left((a_{i;r}^f - 1) \log f_{i;r} - f_{i;r} \frac{a_{i;r}^f}{b_{i;r}^f} - a_{i;r}^f \log \frac{b_{i;r}^f}{a_{i;r}^f} - \log \Gamma(a_{i;r}^f) \right) \\ &\quad + \sum_{j;r} \left((a_{j;r}^h - 1) \log h_{j;r} - h_{j;r} \frac{a_{j;r}^h}{b_{j;r}^h} - a_{j;r}^h \log \frac{b_{j;r}^h}{a_{j;r}^h} - \log \Gamma(a_{j;r}^h) \right) \\ &\quad + \sum_{k;r} \left((a_{k;r}^d - 1) \log d_{k;r} - d_{k;r} \frac{a_{k;r}^d}{b_{k;r}^d} - a_{k;r}^d \log \frac{b_{k;r}^d}{a_{k;r}^d} - \log \Gamma(a_{k;r}^d) \right) \end{aligned} \quad (22)$$

362 4.5.1 Update Rule for $q(\mathcal{W})$

Considering the terms in the log-likelihood expression in Equation 22, that only includes $w_{i;j;k;r}$, we find that

$$\begin{aligned} q(w_{i;j;k;r}) &\propto \exp \left(m_{j;k} \log \left(y_{j;k} - \sum_{i;r} w_{i;j;k;r} \right) \right. \\ &\quad \left. + \sum_{i;r} m_{j;k} \left(w_{i;j;k;r} \log s_{i;j;k;r} - s_{i;j;k;r} - \log \Gamma(w_{i;j;k;r} + 1) \right) \right) \\ &\propto \text{Multinomial}(w_{j;k}; \mathbf{x}_{i;j;k}; \mathbf{p}_{i;j;k;r})^{m_{j;k}} \end{aligned} \quad (23)$$

where, $W_{j;k;r}$ becomes multinomial distributed. The expectation of \mathcal{W} is calculated as

$$\rho_{i;j;k;r} = \frac{\exp(s_{;i} + \langle \log f_{i;r} \rangle + \langle \log h_{j;r} \rangle + \langle \log d_{k;r} \rangle)}{\sum_{i;r} \exp(s_{;i} + \langle \log f_{i;r} \rangle + \langle \log h_{j;r} \rangle + \langle \log d_{k;r} \rangle)} \quad (24)$$

$$\langle W_{i;j;k;r} \rangle = y_{j;k} \rho_{i;j;k;r} \quad (25)$$

363 4.5.2 Update Rule for $q(\mathbf{F})$

Similarly, considering the terms in log-likelihood Equation 22 that only includes $f_{i;r}$, we find that

$$q(f_{i;r}) \propto \left(\sum_{j;k} m_{j;k} \langle W_{i;j;k;r} \rangle + a_{i;r}^f - 1 \right) \log f_{i;r} - \left(\sum_{j;k} m_{j;k} s_{;i} \langle h_{j;r} \rangle \langle d_{k;r} \rangle + \frac{a_{i;r}^f}{b_{i;r}^f} \right) f_{i;r} \quad (26)$$

$$\propto \text{Gamma}(f_{i;r}; \frac{f_{i;r}}{b_{i;r}^f}, \frac{f_{i;r}}{a_{i;r}^f}) \quad (27)$$

where $f_{i;r}$ becomes Gamma distributed with shape and scale parameters

$$\frac{f_{i;r}}{b_{i;r}^f} = a_{i;r}^f + \sum_{j;k} m_{j;k} \langle W_{i;j;k;r} \rangle \quad (28)$$

$$\frac{f_{i;r}}{a_{i;r}^f} = \left(\frac{a_{i;r}^f}{b_{i;r}^f} + \sum_{j;k} m_{j;k} s_{;i} \langle h_{j;r} \rangle \langle d_{k;r} \rangle \right)^{-1} \quad (29)$$

We calculate the expectation of $f_{i;r}$ and the logarithm of $f_{i;r}$ as

$$\langle f_{i;r} \rangle = \frac{f_{i;r}}{b_{i;r}^f} \frac{f_{i;r}}{a_{i;r}^f} \quad (30)$$

$$\langle \log f_{i;r} \rangle = \Psi\left(\frac{f_{i;r}}{b_{i;r}^f}\right) + \log \frac{f_{i;r}}{a_{i;r}^f} \quad (31)$$

364 The variational Bayes algorithm that uses the above equations is given
365 in Algorithm 4. The calculation of the lower bound is given in Appendix 1.
366 The exact derivations of all equations can be found in [28].

367 4.6 Computational Complexity

368 The nonnegative tensor factorization is an NP-hard problem [29]. The varia-
369 tional Bayes algorithm we introduced in Algorithm 4 is an iterative solution
370 that converges to a local maximum solution. The complexity of each itera-
371 tion is determined by the leading term, which is the Equation 25. In general,

Algorithm 4 Variational Bayes Algorithm

```
1: function THINNTF_VB( $\mathcal{Y}, \mathbf{S}, \mathbf{A}^F; \mathbf{B}^F; \mathbf{A}^H; \mathbf{B}^H; \mathbf{A}^D; \mathbf{B}^D$ )
   // Randomly initialize factors and latent tensor
2:    $\{\mathbf{F}; \mathbf{H}; \mathbf{D}; \mathcal{W}\} \leftarrow \text{RANDINIT}(\mathbf{S}; \mathbf{A}^F; \mathbf{B}^F; \mathbf{A}^H; \mathbf{B}^H; \mathbf{A}^D; \mathbf{B}^D)$ 
3:   repeat
4:     Calculate  $\frac{f}{i;r}; \frac{f}{i;r}$  and  $\langle f_{i;r} \rangle$  as in Equations 28, 29 and 30.
5:     Calculate  $\frac{h}{j;r}; \frac{h}{k;r}$  and  $\langle h_{j;r} \rangle$  similarly.
6:     Calculate  $\frac{d}{k;r}; \frac{d}{k;r}$  and  $\langle d_{k;r} \rangle$  similarly.
7:     Calculate  $\langle \log f_{i;r} \rangle$  as in Equation 31.
8:     Calculate  $\langle \log f_{i;r} \rangle$  similarly.
9:     Calculate  $\langle \log f_{i;r} \rangle$  similarly.
10:    Calculate  $\langle W_{:ij;k;r} \rangle$  as in Equation 25.
11:    Calculate lower bound
12:  until Max iterations are reached or lower bound converged
13:  return  $\mathbf{F}, \mathbf{H}, \mathbf{D}, \mathcal{X}$ 
```

372 calculating a ThinNTF model with R components for a d dimensional tensor
373 with all dimensions of length N has $O(N^{(d+1)}R)$ complexity for a single
374 iteration.

375 5 Data Collection

376 Applying a tensor model to the flow length estimation problem requires high
377 volume data collected over a long period, to capture the timely behavior
378 of the network. The already available online data sets do not fulfill this
379 requirement. Therefore, we collected our own real-world data from a mobile
380 network service provider in Turkey [30]. The architecture of our system and
381 the description of the data we collected as presented as follows.

382 5.1 System Architecture

383 The system architecture of a mobile operator's general packet radio service
384 (GPRS) network infrastructure, including radio access and core network el-
385 ements, is illustrated in Figure 7. IP traffic generated or received by mobile
386 devices between mobile station (MS) and packet data network (PDN), e.g.
387 IP Multimedia Subsystem (IMS), is tunneled in the core network of mobile
388 operators through serving GPRS support node (SGSN) and gateway GPRS
389 support node (GGSN) via the user data part of the GPRS tunneling protocol
390 (GTP) [31]. The GTP message exchanges include information such as the

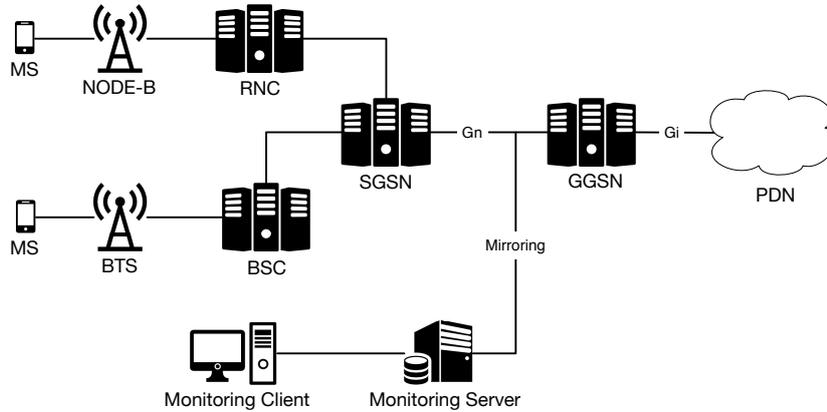


Figure 7: The placement of our monitoring server inside the premises of the mobile operator running a commercial cellular network.

391 size of the traffic, IP session start and end time, device and user identifiers.

392 The Gn interface ¹ carries user packets to be transferred between the
 393 mobile users and the internet together with control packets necessary for the
 394 universal mobile telecommunications service (UMTS) core network [32]. All
 395 packets in this channel are carried by the GTP, which is an IP based protocol
 396 for carrying GPRS data within UMTS networks, used for data encapsulation
 397 in order to keep the core network independent of the protocols that are being
 398 used between MS and the packet-switched network.

399 The Gn interface carries mainly two types of GTP message structures or
 400 packets: GTP-C and GTP-U. GTP-C is used for signaling between SGSN
 401 and GGSN in core network which carries packet data protocol (PDP) Context
 402 messages such as activating and deactivating user session, configuring service
 403 parameters or updating the session. GTP-U is used for transmitting user data
 404 between the radio access network and core network. In our experiments, we
 405 filtered out GTP-C packets (since they are not considered to be part of a flow
 406 due to flow definition), which makes 10% of the total Gn traffic. Therefore,
 407 the sampling is applied to GTP-U packets only. GTP is carried mainly over
 408 UDP.

409 5.2 Data Extraction process

410 The mobile operators network consists of several districts with more than
 411 10 regional core areas through-out Turkey. The average total traffic in all
 412 regional areas consists of approximately over 15 billion packets in the uplink

¹Gn is an interface between SGSN and GGSN where GTP is the main protocol for network packets flowing through.

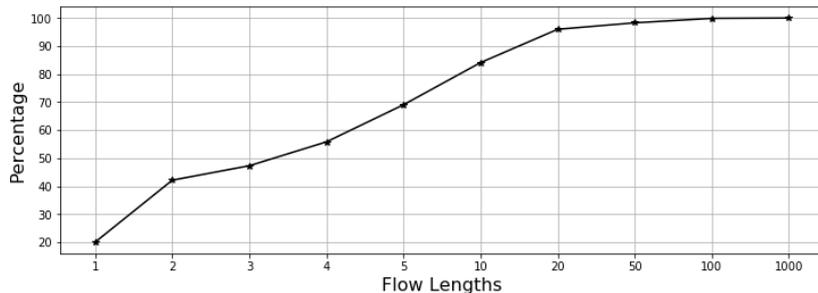


Figure 8: Cumulative flow lengths in the real world data.

413 direction and over 20 billion packets in the downlink direction daily. This
 414 corresponds to approximately 80 terabytes of total data flowing in uplink
 415 and downlink daily inside the mobile operators core network as a whole.
 416 In this work, the Gn interface which connects the SGSN and GGSN nodes
 417 are mirrored and the network traffic is transferred into a FLD server lo-
 418 cated at mobile operator’s technology center in the core network. A speed
 419 of 200Mbit/sec at peak hours can be observed through one of the mirrored
 420 interfaces in the core network.

421 We monitored the network traffic in one of the servers of a mobile operator
 422 continuously for 10 days. We developed a packet extraction tool inside the
 423 monitoring server shown in Figure 7 which parses each GTP-U packet and
 424 stores the packet signature together with packet length and arrival time,
 425 discarding their payload. The stored network data is processed offline for
 426 extracting the true flow lengths.

427 After the data collection and flow extraction, the total number of packets
 428 collected is found to be 4×10^{11} , which makes up around 2.5×10^{10} flows.
 429 Figure 8 shows the cumulative flow length distribution of the data. We see
 430 that most of the flows have less than 5 packets, and 99.9% of them have less
 431 than 100 packets.

432 6 Experiments and Results

433 We designed two sets of experiments in order to verify our model: synthetic
 434 and real-world experiments. In each set, we sampled the original data with
 435 both uniform and ANLS models with different sampling parameters. Then
 436 we tried to recover the original tensor with ThinNTF models. The ThinNTF
 437 model takes a single parameter R which is the number of components in each
 438 factor. Additionally, we also represented data as $I \times JK$ matrix by unfold-
 439 ing the \mathcal{X} tensor in the first dimension as described in [28], and applied the

440 2-dimensional version of ThinNTF, which we simply call thin nonnegative
 441 matrix factorization (ThinNMF). For Uniform sampling, we used the maxi-
 442 mum likelihood estimation (MLE) defined in [3] as the baseline. For ANLS,
 443 we used both MLE and its own unbiased estimator of the model as baselines.

444 Both ThinNMF and ThinNTF models explain the data as a linear com-
 445 bination of R flow length distributions, stored in the columns of \mathbf{F} matrix.
 446 In the ThinNMF model, we have JK coefficient sets for this combination.
 447 On the other hand, in ThinNTF, we have J coefficients for hour-of-day and
 448 K coefficients for day-of-week. The Cartesian product of these coefficient
 449 sets make a total of JK coefficients and creates a dependency between the
 450 hour and day attributes. Therefore, we expect that ThinNTF captures the
 451 periodicity and give better estimates.

452 During the experiments, we always run the stochastic algorithms, i.e.
 453 ThinNMF, ThinNTF, and MLE, for 10 times and keep the parameters of the
 454 model with the highest lower bound value. Then we reported the success of
 455 our algorithm with the weighted mean relative distance (WMRD) metric as
 456 this was used in all previous flow size estimation works. The WMRD is a
 457 metric which gives more weights to the relative differences that occur with
 458 larger frequency. It is formulated as

$$\text{wmrd}(\mathbf{x}; \hat{\mathbf{x}}) = \frac{\sum_i |x_i - \hat{x}_i|}{\sum_i (x_i + \hat{x}_i)} \quad (32)$$

459 where \mathbf{x} is an original flow size distribution measured at the end of the
 460 hour and $\hat{\mathbf{x}}$ is its estimated version. For the whole tensor \mathcal{X} , we calculate the
 461 average WMRD value.

462 Additionally, we report the Kullback-Leibler divergence between the orig-
 463 inal and the estimated tensors, since this is the metric minimized during the
 464 variational Bayes algorithm. The KL divergence between two distributions
 465 \mathbf{x} and $\hat{\mathbf{x}}$ is calculated as

$$\text{KL}(\mathbf{x}; \hat{\mathbf{x}}) = \sum_i x_i \log \left(\frac{\hat{x}_i}{x_i} \right) \quad (33)$$

466 6.1 Experiments on Synthetic Data

467 We prepared our synthetic experiments to test the validity of our models. In
 468 this experiment set, we used the generative model of the ThinNTF model as
 469 described in Algorithm 3 to generate a small network with maximum flow
 470 size $l = 10; J = 24$ and $K = 7$. The original synthetic flow length distribu-
 471 tion \mathcal{X} is generated by a generative model with 3 components, where each

Period	ThinNMF-R3	ThinNTF-R3	MLE
2	0.53	0.49	0.88
4	0.63	0.59	1.20
8	0.65	0.61	1.29
16	0.68	0.61	1.41
32	0.74	0.61	1.37
64	0.85	0.61	1.50

Table 3: Uniform Sampling Results on Synthetic Data

U	ThinNMF-R3	ThinNTF-R3	MLE	ANLS
0.01	0.29	0.27	0.09	0.15
0.02	0.31	0.29	0.17	0.27
0.05	0.33	0.32	0.36	0.28
0.1	0.35	0.34	0.51	0.38
0.2	0.38	0.36	0.59	0.67
0.5	0.48	0.47	0.72	0.70

Table 4: ANLS Sampling Results on Synthetic Data

472 component is a column in the \mathbf{F} factor. We selected these 3 components as
473 exponential, inverted exponential, and uniform random distributions. There-
474 fore, in experiments, we used ThinNMF-R3 and ThinNTF-R3 models, where
475 the suffix R3 shows that the model has 3 components.

476 We sampled the synthetic data with uniform and ANLS sampling meth-
477 ods with different sampling parameters. The sampling was done simply by
478 randomly drawing a sampled size for each flow according to the sampling
479 probabilities in the \mathbf{S} matrix. By this way, we ignored the flow splitting
480 problem and this gave us an ideal data for the ThinNTF model. We report
481 and compare the mean standard deviation of these WMRD values for all
482 experiments.

483 The ThinNTF model always performed best with the uniform sampling
484 model, as shown in Table 3 as expected. On ANLS sampling, the MLE
485 and the ANLS estimators performed better with high sampling probabilities,
486 when $u \in (0.01; 0.02; 0.05)$, as shown in Table 4. On the other hand, when
487 the sampling probability of the ANLS model decreases, the ThinNMF helped
488 with better estimations. From the initial results, we conclude that the factor-
489 ization is definitely helpful for more difficult uniform sampling method and
490 helps lower the sampling probabilities in flow-based packet sampling. The
491 results are also visible in Figure 9.

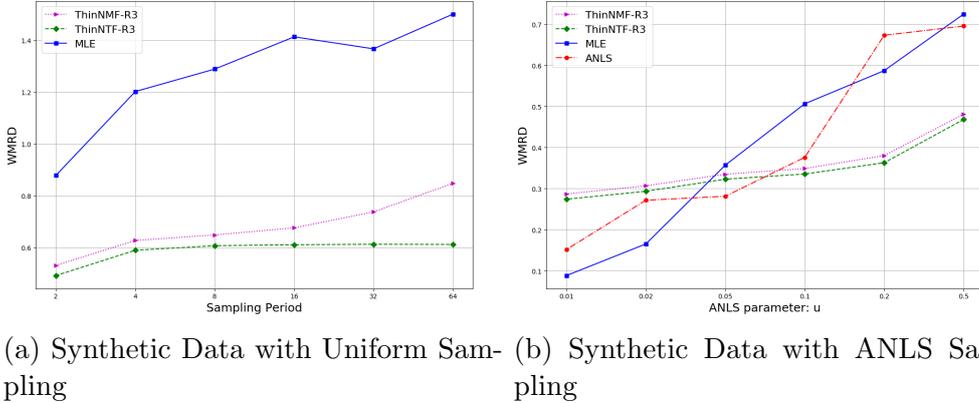


Figure 9: Synthetic experiment results.

492 6.2 Experiments on Real World Data

493 The original data collected from a mobile network provider as we described
 494 in Section 5, is sampled with both sampling methods. However, this time
 495 we simulated the real network offline by feeding the packets one by one to
 496 the monitoring server, as described in Section 5, This way, we were able to
 497 create the actual conditions on a sampler installed at a network provider’s
 498 backbone. This also created the flow splitting problem, since we applied
 499 a 30 seconds time-out in our flow hash. We set the maximum flow length
 500 as $l = 100$, meaning that $\mathcal{X}_{100, \dots}$ entries show the count of flows that have
 501 more than 99 packets. This clamping decision was made according to the
 502 cumulative distribution of flow lengths as shown in Figure 8 We also clamped
 503 the sampling matrices \mathbf{S} so that they exactly match the model.

504 Since the number of components in the original flow distribution is un-
 505 known, we run our experiments with $R \in [2; 3; 4]$ components for ThinNMF
 506 and ThinNTF. The rest of the experiment is similar to the synthetic one.
 507 The sampled Y matrix with shape $100 \times 24 \times 7$ is factorized and $\hat{\mathcal{X}}$ is recon-
 508 structed with the estimated factors. We reported and compared the mean
 509 and standard deviation of 24×7 WMRD and KL values.

510 The factorization models, both ThinNMF and ThinNTF helped lower the
 511 WMRD score in both uniform and ANLS sampling methods. ThinNTF-R4
 512 model consistently gave lower error than the MLE baseline for uniform model
 513 as shown in Table 5 and Figure 10. Indeed, our factorization framework
 514 improved results overall for uniform sampling. However, since recovering
 515 true estimates in uniform sampling is quite difficult, we see less impact of
 516 the factorization as the sampling ratio increases.

517 Figure 10 also gives the KL values between the true and estimated flow

Period	ThinNMF			ThinNTF			MLE
	R=2	R=3	R=4	R=2	R=3	R=4	
2	0.23	0.24	0.23	0.21	0.25	0.22	0.41
4	0.55	0.52	0.53	0.50	0.48	0.49	0.69
8	0.94	0.93	0.94	0.91	0.90	0.87	0.97
16	1.15	1.11	1.11	1.09	1.05	1.04	1.05
32	1.25	1.24	1.24	1.16	1.13	1.10	1.22
64	1.31	1.29	1.30	1.09	1.06	1.04	1.22

Table 5: Uniform Sampling Results on Real Data

U	ThinNMF			ThinNTF			MLE	ANLS
	R=2	R=3	R=4	R=2	R=3	R=4		
0.01	0.03	0.02	0.01	0.05	0.04	0.03	0.05	0.12
0.02	0.04	0.03	0.02	0.06	0.04	0.03	0.08	0.21
0.05	0.04	0.03	0.02	0.07	0.05	0.04	0.13	0.39
0.1	0.06	0.05	0.04	0.08	0.07	0.05	0.17	0.61
0.2	0.08	0.08	0.08	0.10	0.09	0.07	0.21	0.70
0.5	0.13	0.13	0.11	0.16	0.15	0.13	0.33	0.94

Table 6: ANLS Sampling Results on Real Data

length distributions. While the scale of this metric is different, it gives consistent results with the WMRD. This shows that our model, which minimizes the KL metric also minimizes the commonly used WMRD metric, hence the model is suitable for this problem.

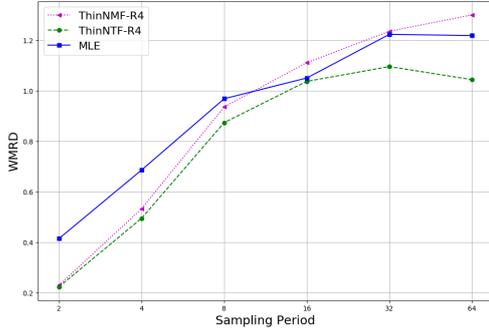
Another important issue is that for uniform sampling, 3-way factorization is more successful than the 2-way factorization. The periodicity information which is captured by the ThinNTF model help improve the estimates and makes it a more successful model for this sampling method.

In ANLS, all our factorization models gave lower error values than the MLE and unbiased estimator of ANLS as shown in Table 6 and Figure 11. Since ANLS is a more powerful sampling method than uniform sampling, our the effect framework is slightly less for small sampling parameter u . However, both ThinNMF and ThinNTF gave better result while sampling smaller number of packets (when u is large). Furthermore, since we are trying to recover the same original data in both experiments, we can compare our ThinNMF and ThinNTF models under two sampling methods. We see that in both methods as the number of components increases, the models gave lower error rates. However, with the uniform sampling method, 3-dimensional methods give better results, while with ANLS, 2-dimensional models perform slightly better.

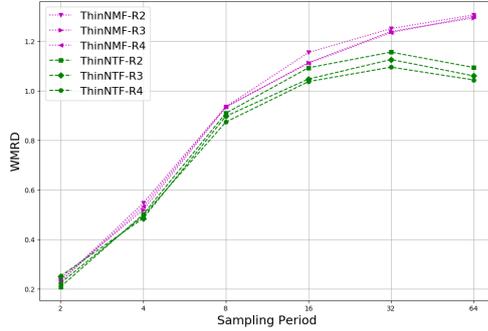
6.3 Effect of Clamping

The choice of where to clamp the data can be given by multiple factors. First of all, one can set the clamping value l_{max} according to a value of special interest. Otherwise, we would like to choose a small l_{max} so that we deal with a dense tensor and we deal with less parameters. On the other hand, we would like to set l_{max} as high as possible so that the clamped portion of the data is as small as possible.

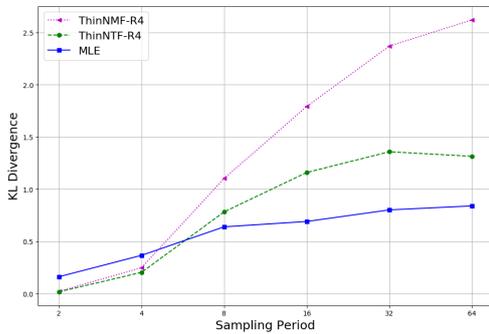
We run the best algorithms found in previous section for uniform and ANLS sampling methods with $l_{max} \in \{25; 50; 75; 100\}$. The WMRD values are given in Figure 12. In both methods, $l_{max} = 25$ gave relatively poor performance and $l_{max} = 100$ was generally the best choice. Also the results with $l_{max} \geq 50$ are closer to each other. This is consistent with the graph in Figure 8, where the cumulative flow lengths do not change much after $l_{max} = 50$. A final remark from this experiment is that as the clamping value increases estimation becomes harder with small sampling rates. This explains the results in uniform sampling with sampling rate $1/64$.



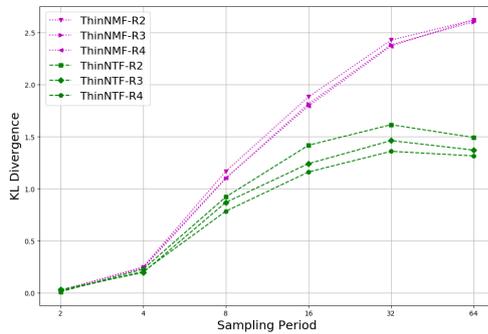
(a) WMRD metrics.



(b) WMRD metrics wrt. R .



(c) KL metrics



(d) KL metrics wrt. R .

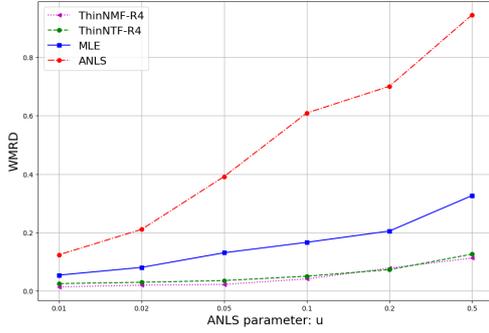
Figure 10: Real world data results with Uniform sampler.

554 7 Conclusions

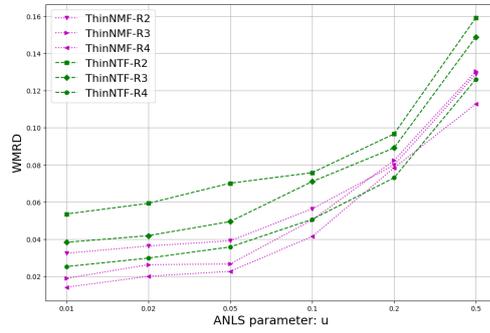
555 In this work, we introduced a novel nonnegative tensor factorization model
 556 called ThinNTF, which extends the classic nonnegative tensor factorization
 557 with an additional constant factor that can represent a network packet sam-
 558 pling method. We showed that this model can be employed to improve the
 559 current reconstruction algorithms in recovering the original flow length dis-
 560 tributions.

561 We tested our model with two different types of sampling methods: the
 562 uniform packet sampling method and a flow-based packet sampling method,
 563 called ANLS. We described how to use these methods by showing how to
 564 build their sampling matrices.

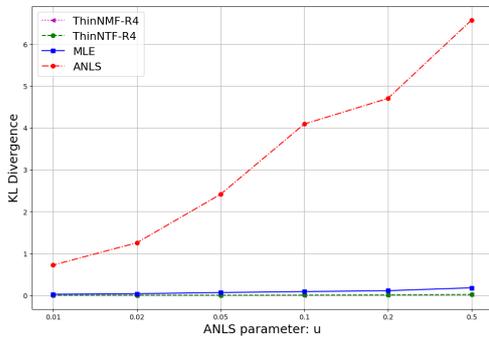
565 In order to test our model, we collected high-volume data from a mobile
 566 network provider for a long period in order to observe the periodical behavior
 567 of the flow length distribution. In experiments on synthetic and real-world
 568 data, our models gave promising results by lowering the estimation errors
 569 compared to the baselines of each sampling method. We conclude that our



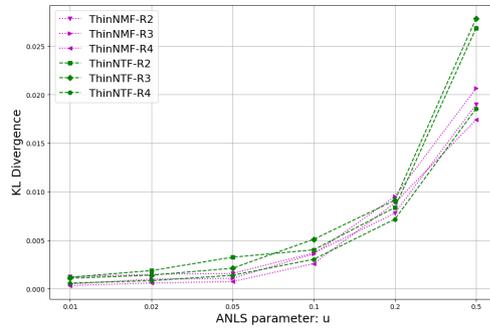
(a) WMRD metrics.



(b) WMRD metrics wrt. R .

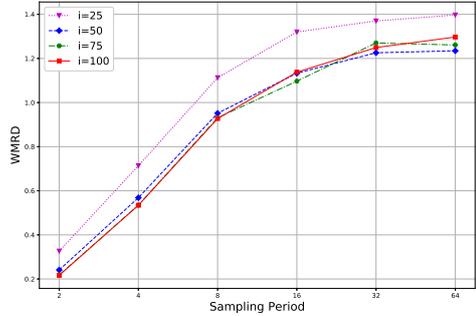


(c) KL metrics

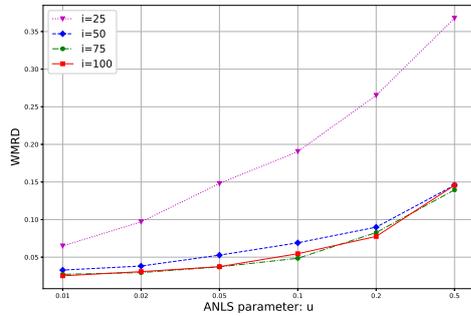


(d) KL metrics wrt. R .

Figure 11: Real world data results with ANLS sampler.



(a) Uniform Sampling ThinNMF-R4



(b) ANLS Sampling ThinNTF-R4

Figure 12: Clamping Experiments.

570 model can be used to decrease estimation errors or to decrease the sampling
 571 probabilities without increasing the estimation error.

572 An important issue left as future work is the online execution of the
 573 ThinNTF model. Theoretically, the ThinNTF model can be used online
 574 once sufficient data from the target network is collected and the flow length

575 distribution components, ie. the \mathbf{F} factor, are inferred. The power of our
576 model is that this inference can be done directly from the sampled obser-
577 vations. Once the \mathbf{F} factor is estimated, for each incoming observation the
578 corresponding entries in other factors can be inferred by keeping \mathbf{F} constant
579 during the inference. Moreover, \mathbf{F} can be updated periodically, say weekly, in
580 a sliding window fashion and kept up to date with the networks flow length
581 behavior.

582 **8 Conflict of Interest**

583 The authors declare that there are no conflicts of interest regarding the pub-
584 lication of this paper.

585 **Appendix A Variational Lower Bound Cal-** 586 **ulation**

587 The calculation of the lower bound includes a few arithmetic tricks. We pro-
588 vide a Bayesian nonnegative matrix factorization [28] tutorial for the detailed
589 derivation and coding tricks. The final form of the lower bound equation is

$$\begin{aligned}
\mathcal{L} &= \langle \log p(\mathcal{Y}; \mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D}) \rangle_{q(\mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D})} + H_{q(\mathcal{W}; \mathbf{F}; \mathbf{H}; \mathbf{D})} & (34) \\
&= - \sum_{i;j:k;r} m_{j;k} (s_{i;j;k;r} \langle f_{i;r} \rangle \langle h_{j;r} \rangle \langle d_{k;r} \rangle) \\
&\quad + \sum_{i;r} \left(- \langle f_{i;r} \rangle \frac{a_{i;r}^f}{b_{i;r}^f} - a_{i;r}^f \log \frac{b_{i;r}^f}{a_{i;r}^f} - \log \Gamma(a_{i;r}^f) \right) \\
&\quad + \sum_{j;r} \left(- \langle h_{j;r} \rangle \frac{a_{j;r}^h}{b_{j;r}^h} - a_{j;r}^h \log \frac{b_{j;r}^h}{a_{j;r}^h} - \log \Gamma(a_{j;r}^h) \right) \\
&\quad + \sum_{k;r} \left(- \langle d_{k;r} \rangle \frac{a_{k;r}^d}{b_{k;r}^d} - a_{k;r}^d \log \frac{b_{k;r}^d}{a_{k;r}^d} - \log \Gamma(a_{k;r}^d) \right) \\
&\quad - \sum_{j;k} m_{j;k} \log \Gamma(y_{j;k} + 1) \\
&\quad - \sum_{i;j:k;r} m_{j;k} \langle w_{i;j;k;r} \rangle \log(p_{i;j;k;r=S;i}) \\
&\quad + \sum_{i;r} \left(\frac{f_{i;r}}{i;r} (\log \frac{f_{i;r}}{i;r} + 1) + \log \Gamma(\frac{f_{i;r}}{i;r}) \right) \\
&\quad + \sum_{j;r} \left(\frac{h_{j;r}}{j;r} (\log \frac{h_{j;r}}{j;r} + 1) + \log \Gamma(\frac{h_{j;r}}{j;r}) \right) \\
&\quad + \sum_{k;r} \left(\frac{d_{k;r}}{k;r} (\log \frac{d_{k;r}}{k;r} + 1) + \log \Gamma(\frac{d_{k;r}}{k;r}) \right) & (35)
\end{aligned}$$

590 References

- 591 [1] G. Varghese and C. Estan, “The measurement manifesto,” *SIGCOMM*
592 *Comput. Commun. Rev.*, vol. 34, no. 1, pp. 9–14, 2004.
- 593 [2] N. Duffield, “Sampling for passive internet measurement: A review,”
594 *Statistical Science*, vol. 19, no. 3, pp. 472–498, 2004.
- 595 [3] N. Duffield, C. Lund, and M. Thorup, “Estimating flow distributions
596 from sampled flow statistics,” *IEEE/ACM Transactions on Networking*,
597 vol. 13, no. 5, pp. 933–946, 2005.
- 598 [4] B. F. Ribeiro, D. F. Towsley, T. Ye, and J. Bolot, “Fisher information
599 of sampled packets: an application to flow size estimation,” in *Internet*
600 *Measurement Conference*, 2006, pp. 15–26.

- 601 [5] N. Hohn and D. Veitch, "Inverting sampled traffic," in *IMC '03: Pro-*
602 *ceedings of the 3rd ACM SIGCOMM conference on Internet measure-*
603 *ment*. ACM Press, 2003, pp. 222–233.
- 604 [6] L. Yang and G. Michailidis, "Sampled based estimation of network traf-
- 605 fic flow characteristics," in *INFOCOM 2007, 26th IEEE International*
606 *Conference on Computer Communications*, May 2007, pp. 1775–1783.
- 607 [7] Cisco netflow. [Online]. Available:
608 [http://www.cisco.com/c/en/us/products/ios-nx-os-software/ios-](http://www.cisco.com/c/en/us/products/ios-nx-os-software/ios-netflow/index.html)
609 [netflow/index.html](http://www.cisco.com/c/en/us/products/ios-nx-os-software/ios-netflow/index.html)
- 610 [8] ntop. [Online]. Available: <http://www.ntop.org/>
- 611 [9] A. Kumar, M. Sung, J. J. Xu, and E. W. Zegura, "A data streaming al-
- 612 gorithm for estimating subpopulation flow size distribution," *SIGMET-*
613 *RICS Perform. Eval. Rev.*, vol. 33, no. 1, pp. 61–72, 2005.
- 614 [10] C. Hu, B. Liu, S. Wang, Y. Cheng, and Y. Chen, "Anls: Adaptive
- 615 non-linear sampling method for accurate flow size measurement," *IEEE*
616 *Transactions on Communications*, vol. 60, no. 3, pp. 789–798, 2012.
- 617 [11] C. Hu, B. Liu, H. Zhao, K. Chen, Y. Chen, Y. Cheng, and H. Wu,
- 618 "Discount counting for fast flow statistics on flow size and flow volume,"
- 619 *IEEE/ACM Transactions on Networking*, vol. 22, no. 3, pp. 970–981,
620 2014.
- 621 [12] A. Kumar and J. J. Xu, "Sketch guided sampling-using on-line esti-
- 622 mates of flow size for adaptive data collection," in *Proc. 2006 IEEE*
623 *INFOCOM*, 2006.
- 624 [13] A. Kumar, J. Xu, and J. Wang, "Space-code bloom filter for efficient
- 625 per-flow traffic measurement," *IEEE J.Sel. A. Commun.*, vol. 24, no. 12,
626 pp. 2327–2339, 2006.
- 627 [14] C. Hu, S. Wang, J. Tian, B. Liu, Y. Cheng, and C. Yan, "Accurate and
- 628 efficient traffic monitoring using adaptive nonlinear sampling method,"
- 629 in *Proc 2008 IEEE INFOCOM*, 2009.
- 630 [15] B. Ermis and A. T. Cemgil, "A Bayesian Tensor Factorization Model
- 631 via Variational Inference for Link Prediction." *ArXiv*, 2014.
- 632 [16] A. Kumar, M. Sung, J. J. Xu, and J. Wang, "Data streaming algo-
- 633 rithms for efficient and accurate estimation of flow size distribution,"
- 634 *SIGMETRICS Perform. Eval. Rev.*, vol. 32, no. 1, pp. 177–188, 2004.

- 635 [17] D. D. Lee and H. S. Seung, “Learning the parts of objects by non-
636 negative matrix factorization.” *Nature*, vol. 401, pp. 788–91, 1999.
- 637 [18] —, “Algorithms for Non-negative Matrix Factorization,” in *Advances*
638 *in Neural Information Processing Systems 13*, no. 1, 2001, pp. 556–562.
- 639 [19] A. T. Cemgil, “Bayesian inference for nonnegative matrix factorisation
640 models.” *Computational intelligence and neuroscience*, vol. 2009, 2009.
- 641 [20] M. W. Berry, M. Browne, A. N. Langville, V. P. Pauca, and R. J.
642 Plemmons, “Algorithms and applications for approximate nonnegative
643 matrix factorization,” *Computational Statistics & Data Analysis*, vol. 52,
644 no. June 2006, pp. 155–173, 2007.
- 645 [21] W. Xu, X. Liu, and Y. Gong, “Document Clustering Based On Non-
646 negative Matrix Factorization,” in *Proceedings of the 26th Annual In-*
647 *ternational ACM SIGIR Conference on Research and Development in*
648 *Informaion Retrieval*, 2003, pp. 267–273.
- 649 [22] F. L. Hitchcock, “The expression of a tensor or a polyadic as a sum of
650 products,” *J. Math. Phys*, vol. 6, no. 1, pp. 164–189, 1927.
- 651 [23] J. Douglas Carroll and J.-J. Chang, “Analysis of individual differences
652 in multidimensional scaling via an n-way generalization of eckart-young
653 decomposition,” *Psychometrika*, vol. 35, pp. 283–319, 01 1970.
- 654 [24] R. A. Harshman, “Foundations of the parafac procedure: Model
655 and conditions for an” explanatory” multi-mode factor analysis,” *UCLA*
656 *Work. Pap. Phon.*, vol. 16, 11 1969.
- 657 [25] R. Bro, “Parafac. tutorial and applications,” *Chemo-*
658 *metrics and Intelligent Laboratory Systems*, vol. 38,
659 no. 2, pp. 149 – 171, 1997. [Online]. Available:
660 <http://www.sciencedirect.com/science/article/pii/S0169743997000324>
- 661 [26] N. Johnson, A. Kemp, and S. Kotz, *Univariate Discrete Distributions*,
662 ser. Wiley Series in Probability and Statistics. Wiley, 2005.
- 663 [27] D. M. Blei, A. Kucukelbir, and J. D. McAuliffe, “Variational inference:
664 A review for statisticians,” *Journal of the American Statistical Associ-*
665 *ation*, vol. 112, no. 518, pp. 859–877, 2017.
- 666 [28] Bayesian nonnegative matrix factorization tutorial. [Online]. Available:
667 <https://github.com/bariskurt/bptf>

- 668 [29] S. Vavasis, “On the complexity of nonnegative matrix factorization,”
669 *SIAM Journal on Optimization*, vol. 20, no. 3, pp. 1364–1377, 2010.
- 670 [30] B. Kurt, E. Zeydan, U. Yabas, I. A. Karatepe, G. K. Kurt, and A. T.
671 Cemgil, “A network monitoring system for high speed network traffic,”
672 in *2016 13th Annual IEEE International Conference on Sensing, Com-*
673 *munication, and Networking (SECON)*, June 2016, pp. 1–3.
- 674 [31] 3GPP, “3gpp evolved packet system (eps); evolved general packet radio
675 service (gprs) tunnelling protocol for control plane (gtpv2-c); stage 3,”
676 3GPP TS 29.274 13.4.0, Tech. Rep., 2015.
- 677 [32] A. Springer and R. Weigel, *The Umts (Universal Mobile Telecom Stan-*
678 *dard) Physical Layer Basics, Standard, and Frontend Matters*. Berlin,
679 Heidelberg: Springer-Verlag, 2002.