

Sampling

1 Introducing

Consider a N -order tensor, \mathcal{X} with size $L_1 \times L_2 \times \dots \times L_N$, the CP decomposition of this tensor is

$$\mathcal{X} = \llbracket \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket = \sum_{r=1}^R \mathbf{a}_r^{(1)} \circ \dots \circ \mathbf{a}_r^{(N)} \quad (1)$$

where

$$\begin{aligned} \mathbf{A}^{(1)} &= [\mathbf{a}_1^{(1)}, \mathbf{a}_2^{(1)}, \dots, \mathbf{a}_r^{(1)}] \in R^{L_1 \times R} \\ \mathbf{A}^{(2)} &= [\mathbf{a}_1^{(2)}, \mathbf{a}_2^{(2)}, \dots, \mathbf{a}_r^{(2)}] \in R^{L_2 \times R} \\ &\vdots \\ \mathbf{A}^{(N)} &= [\mathbf{a}_1^{(N)}, \mathbf{a}_2^{(N)}, \dots, \mathbf{a}_r^{(N)}] \in R^{L_N \times R} \end{aligned}$$

In equation 1, $\mathbf{A}^{(n)}$ is the factor matrix, and the column size R represents the rank of this tensors, which means that \mathcal{X} can be represented by a sum of these R rank one tensors. The vector $\mathbf{a}_r^{(n)}$ is the r -th column of matrix $\mathbf{A}^{(n)}$. In general, let $\mathbf{a}_{*r}^{(n)}$ be the k -th column of $\mathbf{A}^{(n)}$, $\mathbf{a}_{i_n*}^{(n)}$ be the i_n -th row vector of $\mathbf{A}^{(n)}$, and $a_{i_n k}^{(n)}$ be the element of $\mathbf{A}^{(n)}$.

The element in \mathcal{X} satisfies the following equation:

$$x_{\mathbf{i}} = \sum_{r=1}^R a_{i_1 r}^{(1)} \cdot a_{i_2 r}^{(2)} \cdots a_{i_N r}^{(N)} \quad (2)$$

The indexes vector \mathbf{i} is a shorthand for multi-index (i_1, i_2, \dots, i_N) . We propose two methods, which are wedge sampling and diamond sampling, for estimating the maximum elements $x_{\mathbf{i}}$ in tensor when factor matrices $\mathbf{A}^{(n)}, n = 1, 2, \dots, N$ are given. The wedge sampling is so much like the diamond sampling, and the diamond sampling is an improvement of wedge sampling.

2 Diamond Sampling

The diamond sampling is a way to sample the multi-index $\mathbf{i} : (i_1, i_2, \dots, i_N)$ that is proportional to $x_{\mathbf{i}}^2$. We consider a matrix A with size $m \times n$ as a weighted bipartite graph G_A . Meanwhile, in the adjacent matrix of the graph we use 0 (instead of ∞) to represent the two vertices are not adjacent. And the adjacent matrix of G_A is:

$$\begin{pmatrix} \mathbf{0}_{m \times m} & A \\ A^T & \mathbf{0}_{n \times n} \end{pmatrix}$$

Notation	Explanation
\mathcal{A}	tensor
\mathbf{A}	matrix
$\mathbf{A}^{(n)}$	n -th factor matrix of tensor
$\mathbf{a}_{*r}, \mathbf{a}_r$	k -th column of matrix
\mathbf{a}_{i*}	i -th row of matrix
\mathbf{a}	vector
a_{ir}	element of matrix

Table 1: Notation

2.1 Graph representation

In equation 1, we have N factor matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}$. Those N factor matrices will be represented by a weighted $(N+1)$ -partite graph G_T . And we call those $N+1$ partitions as partition $\bar{V}, V_1, V_2, \dots, V_N$, in which partition V_n has L_n vertexes and partition \bar{V} has R vertexes. We call $v_{i_n}^n$ the i_n -th vertex in partition V_n , and \bar{v}_r the r -th vertex in partition \bar{V} . Every two vertexes from different partition classes $V_i, V_j, i, j \in 1, 2, \dots, N$ are not adjacent. A vertex \bar{v}_r in \bar{V} is only adjacent to vertexes v_i^n in V_n when $a_{i_n r}^{(n)}$ is non-zero with the edge weight $a_{i_n r}^{(n)}$. And the adjacent matrix of G_T is :

$$\begin{pmatrix} \mathbf{0}_{R \times R} & \mathbf{A}^{(1)T} & \dots & \mathbf{A}^{(N)T} \\ \mathbf{A}^{(1)} & \mathbf{0}_{L_1 \times L_1} & \dots & \mathbf{0}_{L_1 \times L_N} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{A}^{(N)} & \mathbf{0}_{L_N \times L_1} & \dots & \mathbf{0}_{L_N \times L_N} \end{pmatrix}$$

Under the graph presentation, we define an event $\varepsilon_{i,r',r}$ consisting of three phases:

- phase 1. Pick an edge $e = (v_{i_1}^1, \bar{v}_r)$;
- phase 2. Walk $N-1$ times from the \bar{v}_r to other partitions V_2, \dots, V_N , arrive at $v_{i_2}^2, v_{i_3}^3, \dots, v_{i_N}^N$;
- phase 3. Walk from $v_{i_1}^1$ to \bar{V} and end in \bar{v}'_r .

When the event $\varepsilon_{i,r',r}$ happened, we give a score according to indexes $\mathbf{i} : (i_1, i_2, \dots, i_N)$. We will assign each phase probabilities so that the final score of $\mathbf{i} : (i_1, i_2, \dots, i_N)$ will be a good estimation of $x_{\mathbf{i}}$.

The sampling algorithm is shown in Algorithm 1.

2.2 Probability of picking edge and walking to other partitions

In this part, we introduce the probabilities in each phase.

- Walking with probability (lines 8–10 of Algorithm 1)

In phase 2 we start from a vertex in \bar{V} to V_i , and phase 3 from a vertex in V_1 to \bar{V} . We choose the path(edge) that we walk through according to its weight. That is picking $r \in \{1, 2, \dots, R\}$ with probability $|a_{i_1 r'}^{(1)}| / \|\mathbf{a}_{i_1*}^{(1)}\|_1$ or given r , picking $i_n \in \{1, 2, \dots, L_n\}$ with probability $|a_{i_n r}^{(n)}| / \|\mathbf{a}_{*r}^{(n)}\|_1$.

- Picking an edge (line 6 of Algorithm 1)

When we pick an edge $e = (v_{i_1}^1, \bar{v}_r)$ in phase 1, we are picking the vertexes pair $(v_{i_1}^1, \bar{v}_r)$. Beforehand, we assign each pair a probability while $a_{i_1 k}^{(1)} \neq 0$:

$$p(i_1, r) = |a_{i_1 k}^{(1)}| \|\mathbf{a}_{i_1*}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1 / \|\mathbf{W}\|_1$$

Algorithm 1 Diamond Sampling with factor matrixes

Given factor matrix $\mathbf{A}^{(n)} \in R^{L_n \times R}, n = 1, 2, \dots, N$.

Let s be the number of samples.

```
1: for all  $a_{i_1 r}^{(1)} \neq 0$  do
2:    $w_{i_1 r} \leftarrow |a_{i_1 r}^{(1)}| \|\mathbf{a}_{i_1 *}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1$ 
3: end for
4:  $\mathbf{X} \leftarrow$  all-zeros tensor of size  $L_1 \times L_2 \dots \times L_N$ 
5: for  $\ell = 1, \dots, s$  do
6:   Sample  $(i_1, r)$  with probability  $w_{i_1 r} / \|\mathbf{W}\|_1$ 
7:   for  $n = 2, \dots, N$  do
8:     Sample  $i_n$  with probability  $|a_{i_n r}^{(n)}| / \|\mathbf{a}_{*r}^{(n)}\|_1$ 
9:   end for
10:  Sample  $r'$  with probability  $|a_{i_1 r'}^{(1)}| / \|\mathbf{a}_{i_1 *}^{(1)}\|_1$ 
11:   $x_{i_1, i_2, \dots, i_N} \leftarrow x_{i_1, i_2, \dots, i_N} + \text{sgn}(a_{i_1 r}^{(1)} \cdot a_{i_2 r}^{(2)} \dots a_{i_N r}^{(N)} \cdot a_{i_1 r'}^{(1)}) a_{i_2 r'}^{(2)} \dots a_{i_N r'}^{(N)}$ 
12: end for
```

Where

$$\|\mathbf{W}\|_1 = \sum_{i_1, r} |a_{i_1 r}^{(1)}| \|\mathbf{a}_{i_1 *}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1$$

Then we pick the pair $(v_{i_1}^1, \bar{v}_r)$ according the probability $p(i_1, r)$.

2.3 Scoring samples

Totally, we do s times sampling, and each sample we will get an coordinate $\mathbf{i} = (i_1, i_2, \dots, i_N)$. If this coordinate has not been sampled previously, let the score of this coordiante in the ℓ -th turn be

$$\mathbf{X}_{\mathbf{i}, \ell} = \text{sgn}(a_{i_1 r}^{(1)} \cdot a_{i_2 r}^{(2)} \dots a_{i_N r}^{(N)} \cdot a_{i_1 r'}^{(1)}) a_{i_2 r'}^{(2)} \dots a_{i_N r'}^{(N)},$$

and add $\hat{x}_{\mathbf{i}} = \mathbf{X}_{\mathbf{i}, \ell}$ into a set where we save the scores. Otherwise, increase $\hat{x}_{\mathbf{i}}$ in the set by $\mathbf{X}_{\mathbf{i}, \ell}$. It is shown in [line 11 of Algorithm 1](#). For \mathbf{i} that not be sampled in ℓ -th turn, we can assume that $\mathbf{X}_{\mathbf{i}, \ell} = 0$. So

$$\hat{x}_{\mathbf{i}} = \sum_{\ell} \mathbf{X}_{\mathbf{i}, \ell}$$

In the next part, we will show that $\hat{x}_{\mathbf{i}}$ is a good estimation of $x_{\mathbf{i}}$ in scale.

2.4 Correctness and error bounds

As we defined previously, the event $\varepsilon_{i, r', r}$ is picking a pair $(v_{i_1}^1, \bar{v}_r)$ then pick paths from V_1 to \bar{V} and some addition pathes from \bar{V} to $V_i, i \in 2, \dots, N$. And we assume that $\mathbf{X}_{\mathbf{i}, \ell}, \ell \in \{1, 2, \dots, s\}$ are independent. Under these assumptions, we give two lemmas.

Lemma 2.1. *The expectation of $\hat{x}_{\mathbf{i}}$ equals to $s \cdot x_{\mathbf{i}}^2 / \|\mathbf{W}\|_1$.*

Proof: The final score $\hat{x}_{\mathbf{i}} = \sum_{\ell=1}^s \mathbf{X}_{\mathbf{i}, \ell}$. And

$$\mathbb{E}[\hat{x}_{\mathbf{i}}] = \mathbb{E}[\sum_{\ell=1}^s \mathbf{X}_{\mathbf{i}, \ell}] = s \mathbb{E}[\mathbf{X}_{\mathbf{i}, 1}] \quad (3)$$

The probability of $\varepsilon_{i,r',r}$ is

$$\begin{aligned}
Pr(\varepsilon_{i,r',r}) &= Pr(\text{pick } (v_{i_1}^1, \bar{v}_r)) \cdot Pr(\text{walk to } \bar{v}_r' | \text{given } v_{i_1}^1) \cdot Pr(\text{walk to } v_{i_2}^2 | \text{given } \bar{v}_r) \cdots Pr(\text{walk to } v_{i_N}^N | \text{given } \bar{v}_r) \\
&= \frac{w_{i_1 r}}{\|\mathbf{W}\|_1} \cdot \frac{|a_{i_1 r'}^{(1)}|}{\|\mathbf{a}_{i_1 *}^{(1)}\|_1} \cdot \frac{|a_{i_2 r}^{(2)}|}{\|\mathbf{a}_{* r}^{(2)}\|_1} \cdots \frac{|a_{i_N r}^{(N)}|}{\|\mathbf{a}_{* r}^{(N)}\|_1} \\
&= \frac{|a_{i_1 r}^{(1)}| \|\mathbf{a}_{i_1 *}^{(1)}\|_1 \|\mathbf{a}_{* r}^{(1)}\|_1 \cdots \|\mathbf{a}_{* r}^{(N)}\|_1}{\|\mathbf{W}\|_1} \cdot \frac{|a_{i_2 r'}^{(2)}|}{\|\mathbf{a}_{i_2 *}^{(2)}\|_1} \cdot \frac{|a_{i_1 r}^{(1)}|}{\|\mathbf{a}_{* r}^{(1)}\|_1} \cdots \frac{|a_{i_N r}^{(N)}|}{\|\mathbf{a}_{* r}^{(N)}\|_1} \\
&= \frac{|a_{i_1 r'}^{(1)} a_{i_1 r}^{(1)} \cdots a_{i_N r}^{(N)}|}{\|\mathbf{W}\|_1}
\end{aligned}$$

We get the probability of one walk:

$$Pr(\varepsilon_{i,r',r}) = \frac{|a_{i_1 r'}^{(1)} a_{i_1 r}^{(1)} \cdots a_{i_N r}^{(N)}|}{\|\mathbf{W}\|_1} \quad (4)$$

Using Equation 3 and Equation 4. The expectation

$$\begin{aligned}
\mathbb{E}[x_i/s] &= \mathbb{E}[\mathbf{X}_{i,1}] \\
&= \sum_r \sum_{r'} Pr(\varepsilon_{i,r',r}) \cdot \text{sgn}(a_{i_1 r'}^{(1)} a_{i_1 r}^{(1)} \cdots a_{i_N r}^{(N)}) a_{i_2 r'}^{(2)} \cdots a_{i_N r'}^{(N)} \\
&= \frac{\sum_r \sum_{r'} |a_{i_1 r'}^{(1)} a_{i_1 r}^{(1)} \cdots a_{i_N r}^{(N)}| \text{sgn}(a_{i_1 r'}^{(1)} a_{i_1 r}^{(1)} \cdots a_{i_N r}^{(N)}) a_{i_2 r'}^{(2)} \cdots a_{i_N r'}^{(N)}}{\|\mathbf{W}\|_1} \\
&= \frac{\sum_r \sum_{r'} a_{i_1 r'}^{(1)} \cdot a_{i_2 r}^{(2)} \cdots a_{i_N r}^{(N)} \cdot a_{i_1 r'}^{(1)} a_{i_2 r'}^{(2)} \cdots a_{i_N r'}^{(N)}}{\|\mathbf{W}\|_1} \\
&= \frac{\{\sum_r a_{i_1 r}^{(1)} \cdot a_{i_2 r}^{(2)} \cdots a_{i_N r}^{(N)}\}^2}{\|\mathbf{W}\|_1} \\
&= \frac{x_i^2}{\|\mathbf{W}\|_1}
\end{aligned}$$

□

Lemma 2.2. Fix $\varepsilon > 0$ and error probability $\sigma \in (0, 1)$. Assuming all entries in factor matrices are nonnegative and at most K . If the number of samples

$$s \leq 3K^{N-1} \|\mathbf{W}\|_1 \log(2/\sigma)/(\varepsilon^2 x_i^2),$$

then

$$Pr[|\hat{x}_i - x_i| \geq \varepsilon x_i^2] \leq \sigma$$

Proof: Let

$$y_i = \sum_{\ell} \mathbf{Y}_{i,\ell} = \sum_{\ell} \mathbf{X}_{i,\ell} / K^{N-1}$$

Where $\mathbf{Y}_{i,\ell}$ is in $[0, 1]$ for $\mathbf{X}_{i,\ell}$ is in $[0, K^{N-1}]$ and y_i is a sum of random variables in $[0, 1]$. Applying the Chernoff bound,

$$Pr[y_i \geq \mathbb{E}[y_i]] < \exp(-\varepsilon^2 \mathbb{E}[y_i]/3)$$

By Lemma ??

$$\mathbb{E}[y_i] = \frac{s x_i^2}{K^{N-1} \|\mathbf{W}\|_1}$$

Algorithm 2 Finding top- t largest value

Given factor matrix $\mathbf{A}^{(n)} \in R^{L_n \times R}, n = 1, 2, \dots, N$.

Let s be the number of samples, t' be the budget.

- 1: Sample the score $\hat{x}_{\mathbf{i}}$ using [Algorithm 1](#) and record the coordinates set Ω_s have been sampled.
- 2: Sort the score to extract

$$\Omega_{t'} = \{\mathbf{i} | \hat{x}_{\mathbf{i}} \geq \hat{x}_{\mathbf{i}'}, \forall \mathbf{i}' \in \Omega_s \setminus \Omega_{t'}\}$$

- 3: Compute the actual value $x_{\mathbf{i}}$ of each coordinate in $\Omega_{t'}$
- 4: Sort the actual value to extract

$$\Omega_t = \{\mathbf{i} | x_{\mathbf{i}} \geq x_{\mathbf{i}'}, \forall \mathbf{i}' \in \Omega_{t'} \setminus \Omega_t\}$$

By the choice of s we have $\mathbb{E}[y_{\mathbf{i}}] = (sx_{\mathbf{i}}^2)/(K^{N-1} \|\mathbf{W}\|_1) \leq 3 \log(2/\sigma)$. Then

$$Pr[y_{\mathbf{i}} \geq \mathbb{E}[y_{\mathbf{i}}]] < \sigma/2$$

By the substitution of $y_{\mathbf{i}} = \sum_{\ell} \mathbf{X}_{\mathbf{i},\ell}/K^{N-1} = \hat{x}_{\mathbf{i}}/K^{N-1}$

$$Pr[\hat{x}_{\mathbf{i}} \cdot \|\mathbf{W}\|_1 \geq s \cdot x_{\mathbf{i}}] < \sigma/2$$

Using the Chernoff lower tail bound and identical reasoning. We get

$$Pr[\hat{x}_{\mathbf{i}} \cdot \|\mathbf{W}\|_1 / s \leq (1 - \epsilon)x_{\mathbf{i}}] \leq \sigma/2$$

□

Theorem 2.1. Fix some threshold τ and error probability $\sigma \in (0, 1)$. Assume all entries in factor matrices are nonnegative and at most K . Suppose $s \geq 12K^{N-1} \|\mathbf{W}\|_1 \log(2L_1L_2 \cdots L_N/\sigma)/\tau^2$. Then with probability at least $1 - \sigma$, the following holds for all indexes $\mathbf{i} = (i_1, i_2, \dots, i_N)$ and $\mathbf{i}' = (i'_1, i'_2, \dots, i'_N) : \text{if } x_{\mathbf{i}} > \tau \text{ and } x_{\mathbf{i}'} < \tau/4, \text{ then } \hat{x}_{\mathbf{i}} > \hat{x}_{\mathbf{i}'}$.

proof:

□

2.5 Finding top- t largest value

We use the sampling score $\hat{x}_{\mathbf{i}}$ and coordinate set \mathbf{i} to find the top- t largest value in \mathfrak{X} when given N factor matrices. Let $\Omega_s = \{\mathbf{i}_j | j = 1, 2, \dots, s\}$ be the coordinates have been sampled.

To reduce the computing, a pre-sort is carried out. It sort the scores $\hat{x}_{\mathbf{i}}$ and reserve the top- t' elements. Obviously, t' , which called budget, is always much larger than t .

Then we compute the actual value $x_{\mathbf{i}}$ of the coordinates in the less small set $\Omega_{t'} = \{\mathbf{i} | \hat{x}_{\mathbf{i}} \geq \hat{x}_{\mathbf{i}'}, \forall \mathbf{i}' \in \Omega_s \setminus \Omega_{t'}\}$. And the top- t largest value's coordinates will be $\Omega_t = \{\mathbf{i} | x_{\mathbf{i}} \geq x_{\mathbf{i}'}, \forall \mathbf{i}' \in \Omega_{t'} \setminus \Omega_t\}$.

The reason to do so is that although the score \hat{x} is a good estimation, the variance is much higher in practical. And the budget t' is a tradeoff between accuracy and computation. The algorithm for finding the top- t largest value is shown in [Algorithm 2](#).

2.6 Finding k-NN for query

In recommendation system, it is demanded to find the most k relevant items for a user. On this occasion, there will be a matrix for all users, for consistency, say $\mathbf{A}^{(1)}$ and $N - 1$ matrices $\mathbf{A}^{(n)}, n = 2, \dots, N$ for items. Each row in matrix $\mathbf{A}^{(1)}$ represents a particular user \mathbf{u} . And the tenor with $N - 1$ order consisting a user \mathbf{u} and items matrices is called a ranking tensor(matrix).

$$\mathbf{x}_u = \llbracket \mathbf{u}, \mathbf{A}^{(2)} \dots, \mathbf{A}^{(N)} \rrbracket = \sum_{r=1}^R u_r \cdot \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(N)} \quad (5)$$

It is an exception to reduce one factor matrix into a vector in previous algorithm. However, when we have so many queries, it is time consuming to do sampling one query by one. Obviously, the difference between different queries in sampling processing is the probability for picking vertexes in partition \bar{V} , in other word, or the frequency number (c_1, c_2, \dots, c_R) of each \bar{v}_r that is expected to be sampled.

In the phase 2(walk from a vertex \bar{v}_r in \bar{V} to $V_i, i \in \{2, \dots, N\}$) of the event $\varepsilon_{i,r',r}$, we call the coordinate $\{i_2, \dots, i_N\}$ a sub-path. For effectively implementation, we use these lists $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_R$ to save the sub-paths of sampled query.

2.6.1 Weight assigning and frequency generating

Instead of pick a random pair (\mathbf{u}, \bar{v}_i) , we use the frequency number sequence (c_1, c_2, \dots, c_R) for phase 1. For each $1 \leq r \leq R$, let $u'_r \leftarrow u_r \parallel \mathbf{a}_{*r}^{(1)} \parallel_1 \parallel \mathbf{a}_{*r}^{(2)} \parallel_1 \dots \parallel \mathbf{a}_{*r}^{(N)} \parallel_1$ the same in phase 1. And then choose c_r to have the expected value of $su'_r / \parallel u' \parallel_1$.

$$c_r = \begin{cases} \lfloor su'_r / \parallel u' \parallel_1 \rfloor, & \text{with probability } \lceil su'_r / \parallel u' \parallel_1 \rceil - su'_r / \parallel u' \parallel_1 \\ \lceil su'_r / \parallel u' \parallel_1 \rceil, & \text{with probability } \lfloor su'_r / \parallel u' \parallel_1 \rfloor - su'_r / \parallel u' \parallel_1 \end{cases}$$

2.6.2 Sub-paths

After the choice of c_r , the sampling method repeat c_r times to sample the the rest indexes (i_1, i_2, \dots, i_N) by randomly walking from the \bar{v}_r to other partitions V_2, \dots, V_N that is of the same probability between different queries. So the sub-path (i_1, i_2, \dots, i_N) is used for the next query to save computation. We give the algorithm in [Algorithm 3](#).

3 Implementing details

In section, we introduce some details get the improvement.

3.1 Data Structure

We arrange the factor matrices in specific storage format. we arrange $\mathbf{A}^{(1)}$ to be saved in row major order, and the rest factor matrices in column major order. Also, the sum of each row in $\mathbf{A}^{(1)}$ and the sum of each column in $\mathbf{A}^{(n)}, n = 2, \dots, N$ will be computed previously.

3.2 Picking the edge

We pick the s edge with probability $w_{i_1 r} / \parallel \mathbf{W} \parallel_1$. The weight matrix \mathbf{W} is saved in row major order also and suppose $\boldsymbol{\rho}$ is the vectorization of \mathbf{w} . Since s is much larger than the length of $\boldsymbol{\rho}$, we use the binary search on the cumulative of $\boldsymbol{\rho}$ via [Algorithm 4](#). We use $\mathbf{e} = (e_1, \dots, e_s)$ to record the sampled indexes, and the vertexes pair is $(e_\ell \% L_1, e_\ell \setminus L_1)$. So the sample result is sorted according to r .

3.3 Random walking

In phase 2, when we walk from \bar{v}_r to other partitions V_n . We use the previous result γ_r , which means we need walk totally γ_r times from \bar{v}_r to each partition. Let $\boldsymbol{\rho} = (|a_1^{(n)} r| / \parallel \mathbf{a}_{*r}^{(n)} \parallel_1, \dots, |a_{L_n}^{(n)} r| / \parallel \mathbf{a}_{*r}^{(n)} \parallel_1)$, then sample γ_r times using vose alias's method.

Algorithm 3 Finding k-NN for a query

Given factor matrix $\mathbf{A}^{(n)} \in R^{L_n \times R}$, $n = 1, 2, \dots, N$. $\mathbf{A}^{(1)}$ is the query matrix.

Let s be the number of samples, t' be the budget.

```
1: Initialize  $R$  empty sub-paths lists  $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_R$  for each  $\bar{v}_r$ .
2: for  $i_1 = 1, 2, \dots, L_1$  do
3:   Let query  $\mathbf{u} = \mathbf{a}_{i_1*}^{(1)}$ 
4:   for  $r = 1, \dots, R$  do
5:     Let  $u'_r \leftarrow |u_r| \|\mathbf{a}_{*r}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1$ 
6:   end for
7:   for  $r = 1, \dots, R$  determine  $c_r$  do
8:     
$$c_r = \begin{cases} \lfloor su'_r / \|\mathbf{u}'\|_1 \rfloor, & \text{with probability } \lceil su'_r / \|\mathbf{u}'\|_1 \rceil - su'_r / \|\mathbf{u}'\|_1 \\ \lceil su'_r / \|\mathbf{u}'\|_1 \rceil, & \text{with probability } \lfloor su'_r / \|\mathbf{u}'\|_1 \rfloor - su'_r / \|\mathbf{u}'\|_1 \end{cases}$$

9:   end for
10:  for  $r = 1, \dots, R$  do
11:    if  $c_r \leq |\mathbf{g}_r|$  then
12:      Use  $c_r$  sub-paths in  $\mathbf{g}_r$  as the sampled coordinates.
13:    else
14:      Sample  $n = c_r - |\mathbf{g}_r|$  sub-paths and append it into  $\mathbf{g}_r$ .
15:    end if
16:    Use the method in Algorithm 2 to find the top- $k$  largest value for query  $\mathbf{u}$ .
17:  end for
18: end for
```

Algorithm 4 Picking Edge

Given probability $\mathbf{W} \in R^{L_1 \times R}$, $\boldsymbol{\rho} \in [0, 1]^p$ is the vectorization, and number of samples s . The result saved in \mathbf{e}

```
1: Sample  $\mu_\ell \sim U(0, 1)$  for  $\ell \in \{1, \dots, s\}$ 
2: Sort  $\boldsymbol{\mu}$  so that  $\mu_1 \leq \dots \leq \mu_s$ 
3:  $k \leftarrow 1, \bar{\mu} \leftarrow \mu_k, \boldsymbol{\gamma} \in R^R \leftarrow$  all-zero vector.
4: for  $\ell = 1, \dots, s$  do
5:   while  $\mu_\ell > \bar{\mu}$  do
6:      $k \leftarrow k + 1, \bar{\mu} \leftarrow \mu_k$ 
7:   end while
8:    $r = k \setminus L_1$ 
9:    $e_\ell \leftarrow k, \gamma_r \leftarrow \gamma_r + 1$ 
10: end for
```
