

# 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 an improvement of wedge sampling.

## 2 Diamond Sampling

The diamond sampling is a way to sample the multi-index  $\mathbf{i} : (i_1, i_2, \dots, i_N)$  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  $\infty$  to represent the two vertices are not adjacent. So 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 are represented by a weighted  $(N + 1)$ -partite graph  $G_T$ . And we call those  $N + 1$  partitions as  $\bar{V}, V_1, V_2, \dots, V_N$ . Partition  $V_n$  has  $L_n$  vertexes and partition  $\bar{V}$  has  $R$  vertexes which is the same size of tensor rank. 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  $v_{i_n}^n$  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  $\varepsilon_{i,r',r}$ : 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  $\bar{v}_{r'}$ . After this event happened, we give a score according to indexes  $\mathbf{i} : (i_1, i_2, \dots, i_N)$ . We assign each phase a probability so that the final score of  $\mathbf{i} : (i_1, i_2, \dots, i_N)$  will be a good estimation of  $x_{\mathbf{i}}$ .

## 2.2 Probability of picking edge and walking to other partition

- Walking with probability

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, which is given  $i_1$ , pick  $r \in \{1, 2, \dots, R\}$  with probability  $|a_{i_1 r}^{(1)}| / \|\mathbf{a}_{i_1*}^{(1)}\|_1$  or given  $r$ , pick  $i_n \in \{1, 2, \dots, L_n\}$  with probability  $|a_{i_n r}^{(n)}| / \|\mathbf{a}_{*r}^{(n)}\|_1$ .

- Picking an edge

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

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

Where

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

Then we pick the pair according the probability  $p(i_1, r)$ .

### 2.3 Scoring samples

We do  $s$  time samples, and for 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 in the  $\ell$ -th turn be  $\mathbf{X}_{\mathbf{i},\ell} = \text{sgn}(a_{i_1 r}^{(1)} \cdot a_{i_2 r}^{(2)} \cdots a_{i_N r}^{(N)} \cdot a_{i_1 r'}^{(1)} a_{i_2 r'}^{(2)} \cdots a_{i_N r'}^{(N)})$ , and create a variable  $\hat{x}_{\mathbf{i}} = \mathbf{X}_{\mathbf{i},\ell}$ . Otherwise, increase  $\hat{x}_{\mathbf{i}}$  by  $\mathbf{X}_{\mathbf{i},\ell}$ . For  $\mathbf{i}$  that not be sampled in  $\ell$ -th turn, let  $\mathbf{X}_{\mathbf{i},\ell} = 0$ . So

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

Next, we will show that  $\hat{x}_{\mathbf{i}}$  will be a good estimation of  $x_{\mathbf{i}}$  in scale.

### 2.4 Correctness and error bounds

As we defined previously, the event  $\varepsilon_{\mathbf{i},r',r}$  is picking a pair  $(v_{i_1}^1, \bar{v}_r)$  then pick paths from  $V_{i_1}^1$  in  $V_1$  to  $\bar{V}$  and an addition path from  $\bar{v}_r$  in  $\bar{V}$  to  $V_1$ . And we think  $\mathbf{X}_{\mathbf{i},\ell}, \ell \in \{1, 2, \dots, s\}$  are independent. Under this assumption, we give two lemmas.

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

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

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

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

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

We get the probability of one walk.

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

Using equation in 3 and 4.

$$\begin{aligned}
\mathbb{E}[x_{\mathbf{i}}/s] &= \mathbb{E}[\mathbf{X}_{\mathbf{i},1}] \\
&= \sum_r \sum_{r'} Pr(\varepsilon_{\mathbf{i},r',r}) \cdot \text{sgn}(a_{i_1r}^{(1)} \cdot a_{i_2r}^{(2)} \cdots a_{i_Nr}^{(N)} \cdot a_{i_1r'}^{(1)} a_{i_2r'}^{(2)} \cdots a_{i_Nr'}^{(N)}) \\
&= \frac{1}{\|\mathbf{W}\|_1} \sum_r \sum_{r'} |a_{i_1r}^{(1)}| \cdot |a_{i_2r}^{(2)}| \cdots |a_{i_Nr}^{(N)}| \cdot |a_{i_1r'}^{(1)}| \text{sgn}(a_{i_1r}^{(1)} \cdot a_{i_2r}^{(2)} \cdots a_{i_Nr}^{(N)} \cdot a_{i_1r'}^{(1)} a_{i_2r'}^{(2)} \cdots a_{i_Nr'}^{(N)}) \\
&= \frac{1}{\|\mathbf{W}\|_1} \sum_r \sum_{r'} a_{i_1r}^{(1)} \cdot a_{i_2r}^{(2)} \cdots a_{i_Nr}^{(N)} \cdot a_{i_1r'}^{(1)} a_{i_2r'}^{(2)} \cdots a_{i_Nr'}^{(N)} \\
&= \frac{1}{\|\mathbf{W}\|_1} \left\{ \sum_r a_{i_1r}^{(1)} \cdot a_{i_2r}^{(2)} \cdots a_{i_Nr}^{(N)} \right\}^2 \\
&= \frac{x_{\mathbf{i}}^2}{\|\mathbf{W}\|_1}
\end{aligned}$$

□

**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

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

then

$$Pr[|\hat{x}_{\mathbf{i}} \|\mathbf{W}\|_1 / s - x_{\mathbf{i}}^2| \geq \varepsilon x_{\mathbf{i}}^2] \leq \sigma$$

**Theorem 2.1.** Fix some threshold  $\tau$  and error probability  $\sigma \in (0, 1)$ . Assume all entries in factor matrices are nonnegative and at most  $K$ . Suppose  $s \geq 12K^N \|\mathbf{W}\|_1 \log(2L_1 L_2 \cdots L_N/\sigma)/\tau^2$ . Then with probability at least  $1 - \sigma$ , the following holds for all indexes  $\mathbf{i} = (i_1, i_2, \dots, i_N)$  and  $\mathbf{i}' = (i'_1, i'_2, \dots, i'_N) : \text{if } x_{\mathbf{i}} > \tau \text{ and } x_{\mathbf{i}'} < \tau/4, \text{ then } \hat{x}_{\mathbf{i}} > \hat{x}_{\mathbf{i}'}.$

## 2.5 Finding top-t largest value

## 2.6 Finding k-NN for query

At this situation, we use on row of matrix  $\mathbf{A}^{(1)}$  each time called  $\mathbf{u}$ . And the tensor of N-1 order called rank tensor for query  $\mathbf{u}$ .

$$\mathbf{x}_{\mathbf{u}} = \left[ \mathbf{u}, \mathbf{A}^{(2)} \dots, \mathbf{A}^{(N)} \right] = \sum_{r=1}^R u_r \cdot \mathbf{a}_r^{(2)} \circ \dots \circ \mathbf{a}_r^{(N)} \quad (5)$$

And find the k most relevant vector sets, each vector set consist N-1 vectors and the

The biggest difference between different queries in sampling processing is the probability for picking vertexes in partition  $\bar{V}$ , or the frequency sequence  $(c_1, c_2, \dots, c_R)$  of  $\bar{v}_r$  to be sampled. Notice that,  $\sum_{r=1}^R c_r = s$ .

For effectively implementation, we use the lists  $\mathbf{g}_1, \mathbf{g}_2, \dots, \mathbf{g}_R$  to record the sub-path(walk from a vertex  $\bar{v}_r$  in  $\bar{V}$  to  $V_i, i \in \{2, \dots, N\}$ ).

### 2.6.1 Weight assigning and frequency generating

For each  $1 \leq r \leq R$ , let  $u'_r \leftarrow u_r \|\mathbf{a}_{*r}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1$ . Instead of sampling  $r$  directly  $s$  times, choosing the  $c_r$  randomly such that  $c_r$  has the expected value  $su'_r / \|\mathbf{u}'\|_1$  will still works.

$$c_r = \begin{cases} \lfloor su'_r / \|\mathbf{u}'\|_1 \rfloor, & \text{with probability } \lceil su'_r / \|\mathbf{u}'\|_1 \rceil - su'_r / \|\mathbf{u}'\|_1 \\ \lceil su'_r / \|\mathbf{u}'\|_1 \rceil, & \text{with probability } su'_r / \|\mathbf{u}'\|_1 - \lfloor su'_r / \|\mathbf{u}'\|_1 \rfloor \end{cases}$$

### 2.6.2 Sub-walks

When given some  $r$ , sampling method need to pick the left indexes  $\mathbf{i}' = (i'_1, i'_2, \dots, i'_N)$ , which has been stored in the previous query, it saves a lot of computation.

- for  $r = 1, 2, \dots, R$ :  
repeat  $c_r$  times: sample  $r'$  times with the probability  $|u_r| / \|u\|_1$ . If the size of  $\mathbf{g}_r$  is larger than  $c_r$ , use the sub-walk in  $\mathbf{g}_r$  as the remaining result. Otherwise, sample the remaining.
- for  $r = 1, 2, \dots, R$ :
- 

## 3 Two Approaches

In this section, we introduce wedge sampling and the advanced edition diamond sampling.

### 3.1 Wedge Approach

A random-sampling based algorithm that identifies high-valued entries of nonnegative matrix products, by assigning scores to each entry of  $\mathbf{q}^T \mathbf{A}$ , where  $\mathbf{A} \in R^{d \times n} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$  is a matrix whose rows are the instance, and  $d$  is the dimension of feature vector. Our task is to find the most relevant instance  $\mathbf{a}_i$  that maximize dot product  $\mathbf{q}^T \mathbf{a}_i$ ,

### 3.2 Graph representation

Consider a product  $\mathbf{M} = \mathbf{Q}\mathbf{A}$ , where  $\mathbf{Q} \in R^{m \times d}, \mathbf{A} \in R^{d \times n}$  are represented by a layered graph with three layers, and treat matrix  $\mathbf{Q}, \mathbf{A}$  as weighted directed bipartite graph. There are  $m$  nodes  $\{v_1^1, v_2^1, \dots, v_m^1\}$  in first layer,  $d$  nodes  $\{v_1^2, v_2^2, \dots, v_d^2\}$  in second layer, and  $n$  nodes  $\{v_1^3, v_2^3, \dots, v_n^3\}$  in the third. For a query vector  $\mathbf{q}^T$ , see it as the  $i$ -th row of  $\mathbf{Q}$ , and  $\mathbf{m} = \mathbf{q}^T \mathbf{A}$ , then define a trace: random walks from the  $i$ -th node, which is the query vector  $\mathbf{q}^T$  in first layer, and terminates at the  $j$ -th node in third layer. In this situation, the first layer has one node.

Walk  $S$  times randomly will terminate in third layer  $S$  times, and record the time node  $j \in \{1, 2, \dots, n\}$  has been reached to  $S_j$ , clearly  $\sum_{j=1}^n S_j = S$ . If  $j \in \{1, 2, \dots, n\}$  is select with probability  $m_j / \sum_{h=1}^n m_h$ , then the expected value of  $S_j$  equals  $m_j S / \sum_{h=1}^n m_h$ . Hence,  $S_j \sum_{h=1}^n m_h / S$  is an unbiased estimator for  $m_j$ .

For every pair of nodes  $\{v_k^2, v_j^3\}$ , assign the probability to the edge

$$p(e_{kj}) = a_{kj} / \sum_{h=1}^d a_{kh}$$

Processing goes in [Algorithm 1](#). Think the trace as first walk to node  $v_k^2$ , and then  $v_j^3$ . So firstly, compute the times  $c_k$  expected to go through node  $v_k^2$ , and sample  $c_k$  times to arrive  $v_j^3$ .

- The expect of  $c_k$ .

$$E[c_k] = \frac{S q'_k}{\sum q'_k}$$

---

**Algorithm 1** Wedge Sampling

---

Given matrix  $\mathbf{A}$  and query  $\mathbf{q}$

Let  $S$  be the number of samples.

```

1: Preprocessing stage
2: for  $k=1, \dots, d$  do
3:   compute  $a_k = \sum_{h=1}^n a_{kh}$  and  $q'_k = q_k a_k$ 
4: end for
5: for  $k=1, \dots, d$  do determine the number of samples  $c_k$  needed for
6:

$$c_k = \begin{cases} \lfloor \frac{Sq'_k}{\sum q'_k} \rfloor, & \text{w/prob. } \lceil \frac{Sq'_k}{\sum q'_k} \rceil - \frac{Sq'_k}{\sum q'_k} \\ \lceil \frac{Sq'_k}{\sum q'_k} \rceil, & \text{w/prob. } \lfloor \frac{Sq'_k}{\sum q'_k} \rfloor - \frac{Sq'_k}{\sum q'_k} \end{cases}$$

7: end for
8: for  $k = 1, \dots, d$  do
9:   Sample  $j \in \{1, \dots, n\}$   $c_k$  times with probability  $p(e_{kj}) = a_{kj} / \sum_{h=1}^n a_{kh}$ .
10:  Increase the counter of  $S_j$ .
11: end for
12: Postprocessing

```

---

- The expect of  $S_j$ .

$$\begin{aligned}
E[S_j] &= \sum_{k=1}^d c_k p(e_{kj}) \\
&= \sum_{k=1}^d \frac{Sq_k a_k}{\sum q_k a_k} \frac{a_{kj}}{a_k} \\
&= \frac{S}{\sum q_k a_k} \sum_{k=1}^d q_k a_{kj} \\
&= \frac{Sm_j}{\sum q_k a_k} = \frac{Sm_j}{\sum m_h}
\end{aligned}$$

As we can see,  $S_j \sum_{h=1}^n m_h / S$  is an unbiased estimator for  $m_j$ .

### 3.3 Implementing details for the wedge sampling

To find the maximum of  $\mathbf{M} = \mathbf{QA}$ , see  $\mathbf{Q}$  as a bunch of query vector  $\mathbf{q}_i$ . Replace the matrix  $\mathbf{A}$  by lists  $L_1, L_2, \dots, L_d$  of indices  $\{1, 2, \dots, n\}$ . For each input query  $\mathbf{q}$ , the algorithm first processes  $\mathbf{q}$  to determine how many samples are needed for each of the  $d$  lists  $L_k$ , or the times  $c_k$  expected to go through node  $v_k^2$ . It then retrieves these samples from the appropriate stored lists and counts how many times each of  $\{1, 2, \dots, n\}$  was sampled. It saves the time for repeating sampling of different query in [line 9 of Algorithm 1](#).

### 3.4 Diamond Approach

For  $\mathbf{m} = \mathbf{q}^T \mathbf{A}$ , the diamond sampling aim to find two intersecting wedges, so that the probability to terminate at  $j = \{1, 2, \dots, n\}$  is proportional to  $m_j^2$ .

## 4 Sampling in Tensor

We can consider the problem as finding some pair of vectors that from two different bunch of vector set, in previous, to have the the top-t maximum dot products.

Now, we extend the problem from two vector set to  $N$  vector sets. And find a list of vectors  $\{\mathbf{a}_{i_1}^{(1)}, \dots, \mathbf{a}_{i_N}^{(N)}\}$ , where  $\mathbf{a}_{i_n}^{(n)}$  is a vector from the  $n$ -th vector set  $\mathbf{A}^{(n)}$  and  $i_n \in \{1, \dots, L^{(n)}\}$ , according to the value of

$$\sum_{k=1}^R a_{i_1 k}^{(1)} a_{i_2 k}^{(2)} \cdots a_{i_N k}^{(N)}$$

A more mathematic interpretation is: given factor matrixes  $\mathbf{A}^{(n)} \in R^{L^{(n)} \times R}$ ,  $n = \{1, 2, \dots, N\}$ , which is the CP decomposition of tensor  $\mathbf{X}$ . Find the top-t value in tensor  $\mathbf{X}$ .

### 4.1 Wedge Sampling in Tensor

We separated those matrixes into three layers, the first layer has  $L^1$  nodes, indicated by the first factor matrix  $\mathbf{A}^{(1)}$ . And the second layer called the middle layer has  $R$  nodes, which picture the feature dimension or the rank of CP decomposition. The layer has  $N - 1$  parallel layers, each layer has  $L^n$  nodes. Under this assumption, we define a trace: walk from the first layers at node  $v_{i_1}^1$  to middle layer  $v_k^m$ , then terminates at each parallel layers with nodes  $\{v_{i_2}^2, v_{i_3}^3, \dots, v_{i_N}^N\}$ .

As similar, define the probability of  $\{v_k^m, v_{i_n}^n\}$ :

$$P(e_{ki_n}) = a_{i_n k}^{(n)} / \sum_{h=1}^R a_{i_n h}^{(n)}$$

#### 4.1.1 Find the most relevant vector lists of query $q$

Given a query  $q$ , find the most relevant vector lists. It is the same when  $\mathbf{A}^{(1)}$  is a vector or see  $q$  as one column of  $\mathbf{A}^{(1)}$ , since we define the row of factor matrix is the number of instances.

- Preprocessing  
 $q'_k = q_k \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1$ , where  $\|\mathbf{a}_{*r}^{(n)}\|_1 = \sum_{h=1}^R a_{i_n h}^{(n)}$ .
- Compute the times  $c_k$  expected to go through node  $v_k^m$

$$c_i = \begin{cases} \lfloor \frac{Sq'_k}{\sum q'_k} \rfloor, & \text{w/prob. } \lceil \frac{Sq'_k}{\sum q'_k} \rceil - \frac{Sq'_k}{\sum q'_k} \\ \lceil \frac{Sq'_k}{\sum q'_k} \rceil, & \text{w/prob. } \lfloor \frac{Sq'_k}{\sum q'_k} \rfloor - \frac{Sq'_k}{\sum q'_k} \end{cases}$$

As we can see

$$E[c_k] = \frac{Sq'_k}{\sum q'_k}$$

- Sample each node in parallel layers  $c_k$  times with probability  $P(e_{ki_n})$ , and increase the counter of  $S_{i_2, \dots, i_N}$

The expect of  $S_{i_2, \dots, i_N}$ .

$$\begin{aligned}
E[S_{i_2, \dots, i_N}] &= \sum_{k=1}^R c_k p(e_{ki_2}) p(e_{ki_3}) \cdots p(e_{ki_N}) \\
&= \frac{S}{\sum q'_k} \sum_{k=1}^R q_k \cdot a_{i_2 k}^{(2)} \cdots a_{i_N k}^{(N)} \\
&= \frac{S m_{i_2, \dots, i_N}}{\sum q'_k} \\
&= \frac{S m_{i_2, \dots, i_N}}{\sum m_{i_2, \dots, i_N}}
\end{aligned}$$

So the  $S_{i_2, \dots, i_N} \sum m_{i_2, \dots, i_N} / S$  is an unbiased estimator for  $m_{i_2, \dots, i_N}$ .

#### 4.1.2 Find the maximum value in tensor

The counter can only indicate the relatively relation of one vector list to query, which can not be used to compare different query. So we need to consider the value overall. Query a single vector, however, is still useful in some application like recommendation system, for at most time we only need to find the most relevant item for on user(query), and the global maximum is useless and time consuming. The trick of replace data instance by a list of indices used in pervious is also helpful when we consider a number of user at a time.

To find the maximum is tensor  $\mathfrak{X}$ , the diamond sampling supply referenced method.

- **Preprocessing**

For every  $i_1, k$  compute the weight

$$w_{i_1 k} = |a_{i_1 k}^{(1)}| \| \mathbf{a}_{*r}^{(2)} \|_1 \cdots \| \mathbf{a}_{*r}^{(N)} \|_1, \text{ where } \| \mathbf{a}_{*r}^{(n)} \|_1 = \sum_{h=1}^R a_{i_n h}^{(n)}.$$

- **Random walk to  $v_k^m$ .**

Sample  $S$  times with the probability  $w_{i_1 k} / \| \mathbf{W} \|_1$

- Sample each node in parallel layers

Given  $k$  sample  $i_2, \dots, i_N$  with probability  $P(e_{ki_n})$ , and increase the value of  $c_{i_1, i_2, \dots, i_N}$  by

$$a_{i_2 k}^{(2)} \cdots a_{i_N k}^{(N)}$$

The expect of  $x_{i_1, i_2, \dots, i_N}$ .

$$\begin{aligned}
E[x_{i_1, i_2, \dots, i_N}] &= \sum_{k=1}^R p(i_1, k) p(e_{ki_2}) p(e_{ki_3}) \cdots p(e_{ki_N}) \\
&= \frac{S}{\| \mathbf{W} \|_1} \sum_{k=1}^R a_{i_1 k}^{(1)} \cdot a_{i_2 k}^{(2)} \cdots a_{i_N k}^{(N)} \\
&= \frac{S}{\| \mathbf{W} \|_1} x_{i_1, i_2, \dots, i_N}
\end{aligned}$$

So the  $c_{i_1, i_2, \dots, i_N} \frac{S}{\| \mathbf{W} \|_1}$  is an unbiased estimator for  $x_{i_1, i_2, \dots, i_N}$ .

## 4.2 Diamond Sampling in Tensor

Diamond sampling is the advanced edition of wedge sampling, which can deal with the negative value.

#### 4.2.1 Find the maximum value in tensor

A direct extension for diamond sampling in tensor is stated in [Algorithm 2](#).



---

**Algorithm 2** Diamond Sampling with factor matrixes

---

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

Let  $s$  be the number of samples.

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

---

#### 4.2.2 Find the most relevant vector lists of query $u$

Consider query  $u$  is one column of  $\mathbf{A}^{(1)}$ , and the retrieve tensor

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

We can use the counter to indicate the value, but is will not deal with the negative value, so a degenerate method of diamond sampling was proposed. The way to find the maximum of  $x_{u,i}$  may be more concise when sampling starts with  $u$  as we can see in [line 2 of Algorithm 3](#).

---

**Algorithm 3** Diamond Sampling with a query vector

---

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

Let  $s$  be the number of samples.

```
1: for all  $u_k \neq 0$  do
2:    $w_k \leftarrow |u_k| \| u \|_1 \| \mathbf{a}_{*r}^{(2)} \|_1 \| \mathbf{a}_{*r}^{(3)} \|_1 \dots \| \mathbf{a}_{*r}^{(N)} \|_1$ 
3: end for
4:  $\mathcal{X}_u \leftarrow$  all-zeros tensor of size  $L^{(2)} \times L^{(3)} \dots \times L^{(N)}$ 
5: for  $\ell = 1, \dots, s$  do
6:   Sample  $k$  with probability  $w_k / \| u \|_1$ 
7:   for  $n = 2, \dots, N$  do
8:     Sample  $i_n$  with probability  $|a_{i_n k}^{(n)}| / \| \mathbf{a}_{*r}^{(n)} \|_1$ 
9:   end for
10:  Sample  $k'$  with probability  $|u_{k'}| / \| u \|_1$ 
11:   $x_{i_2, i_3, \dots, i_N} \leftarrow x_{i_2, i_3, \dots, i_N} + \text{sgn}(a_{i_2 k}^{(2)} \cdot a_{i_3 k}^{(3)} \dots a_{i_N k}^{(N)} \cdot u_k \cdot u_{k'}) a_{i_2 k'}^{(2)} a_{i_3 k'}^{(3)} \dots a_{i_N k'}^{(N)}$ 
12: end for
13: Postprocessing
```

---

## 5 Implementation issues

\*\*\*\*\*

### 5.1 Determined Sample Numbers Online

Pre-sample a number of instance, and it can be used with the accuracy bound to determine a reasonable number of samples.

### 5.2 Dealing with query

When we want to find the maximus relevant vectors for a number of queries. It is effective to save the the sampled instances into lists  $L_1, L_2, \dots, L_d$ . When processes a query, sampled the nodes in middle layer. It then retrieves these samples form the stored lists.

## 6 Experiments

---

\*\*\*\*\*

---

## 7 Sampling in Fashion Recommendation

### 7.1 Structure of $U$ and $A^{(n)}$

Using sampling to solve a fashion recommendation problem.

$$r_{u, \Theta_t} = \sum_{n=1}^N \mathbf{u}^{(n)} \cdot \mathbf{h}^{(n,n)}(\mathbf{x}_{i_n}) + \sum_{n=1}^N \sum_{m=n+1}^N \mathbf{h}^{(n,m)}(\mathbf{x}_{i_n}) \cdot \mathbf{h}^{(m,n)}(\mathbf{x}_{i_m})$$

Given a user  $\mathbf{u}^{(n)} \in R^D$ , and the latent space vector  $\mathbf{h}^{(n,m)}(\mathbf{x}_{i_n}) \in R^D$ , we need to find the maximum value in tensor  $r_{u, \Theta_t}$ , which is the outfit most be like by this user. To do so, we need to reform the vectors into factor matrixes before hand.

We introduce a matrix called  $Matrix(\mathbf{a}_{i_n*}^{(n)})$  with size  $N \times ND$ , which is the reshaped matrix of vector  $\mathbf{a}_{i_n*}^{(n)} \in R^{N^2D}$ .

$$Matrix(\mathbf{a}_{i_n*}^{(n)}) = \begin{matrix} & & & i_n\text{-th column} \\ & & & \frac{1}{2}\mathbf{h}^{(n,1)}(\mathbf{x}_{i_n}) & \cdots & \mathbf{e} \\ & & & \frac{1}{2}\mathbf{h}^{(n,2)}(\mathbf{x}_{i_n}) & \cdots & \mathbf{e} \\ & & & \vdots & \vdots & \vdots \\ i_n\text{-th row} & \frac{1}{2}\mathbf{h}^{(n,1)}(\mathbf{x}_{i_n}) & \frac{1}{2}\mathbf{h}^{(n,2)}(\mathbf{x}_{i_n}) & \cdots & \mathbf{h}^{(n,n)}(\mathbf{x}_{i_n}) & \cdots & \frac{1}{2}\mathbf{h}^{(n,N)}(\mathbf{x}_{i_n}) \\ & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ & \mathbf{e} & \mathbf{e} & \cdots & \frac{1}{2}\mathbf{h}^{(n,N)}(\mathbf{x}_{i_n}) & \cdots & \mathbf{e} \end{matrix} \quad (6)$$

Where  $\mathbf{x}_{i_n}$  is the  $i_n$ -th item in  $n$ -th category,  $\mathbf{h}^{(n,m)}(\mathbf{x}_{i_n})$ , a row vector, is the function to map the  $n$ -th category's items to the latent space for matching with the  $m$ -th category's items,  $\mathbf{e}$  a row vector with all ones. And every vector can be seen as a block of this matrix. Use  $\mathbf{a}_{i_n*}^{(n)}$  to form the factor matrix  $\mathbf{A}^{(n)} \in R^{L^{(n)} \times N^2 D}$ .

And similarly, we define the  $Matrix(\mathbf{U})$ , as the reshaped matrix of row vector  $\mathbf{u}$ .

$$Matrix(\mathbf{u}) = \begin{pmatrix} \mathbf{u}^{(1)} & \mathbf{e} & \cdots & \mathbf{e} \\ \mathbf{e} & \mathbf{u}^{(2)} & \cdots & \mathbf{e} \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{e} & \mathbf{e} & \cdots & \mathbf{u}^{(N)} \end{pmatrix} \quad (7)$$

With the definition of equation 6 and equation 7, we can conclude that:

$$\begin{aligned} r_{u, \Theta_t} &= \sum_{n=1}^N \mathbf{u}^{(n)} \cdot \mathbf{h}^{(n,n)}(\mathbf{x}_{i_n}) + \\ &\quad \sum_{n=1}^N \sum_{m=n+1}^N \mathbf{h}^{(n,m)}(\mathbf{x}_{i_n}) \cdot \mathbf{h}^{(m,n)}(\mathbf{x}_{i_m}) \\ &= \sum_{k=1}^{N^2 D} u_k a_{i_1 k}^{(1)} a_{i_2 k}^{(2)} \cdots a_{i_N k}^{(N)} \\ &= x_i \end{aligned}$$

Suppose  $k-1 = xND + yD + z$ , then  $(x+1, y+1)$  represents the block vector's position in  $Matrix(\mathbf{a}_{i_n*}^{(n)})$  or  $Matrix(\mathbf{u})$ , and  $z+1$  represents the element's index in  $(x+1, y+1)$ -th block. i.e

$$a_{i_n k}^{(n)} = \begin{cases} 1, & x+1 \neq n, y+1 \neq n; \\ (1/2)h_{z+1}^{(n,y+1)}(\mathbf{x}_{i_n}), & x+1 = n, y+1 \neq n; \\ (1/2)h_{z+1}^{(n,x+1)}(\mathbf{x}_{i_n}), & x+1 \neq n, y+1 = n; \\ h_{z+1}^{(n,n)}(\mathbf{x}_{i_n}), & x+1 = y+1 = n. \end{cases} \quad u_k = \begin{cases} 1, & x+1 \neq n \text{ or } y+1 \neq n; \\ u_{z+1}^{(n)}, & x+1 = y+1 = n. \end{cases}$$

Base on the special structure of  $\mathbf{u}$  and  $\mathbf{A}^{(n)}$ , the Algorithm 3 can have more concise formation.

## 7.2 Compute the weight $u_k$

See the line 2 of Algorithm 3

- $\|\mathbf{U}\|_1$  is redundant.
- Only two of  $\{ |u_k|, \|\mathbf{a}_{*r}^{(1)}\|_1, \|\mathbf{a}_{*r}^{(2)}\|_1, \dots, \|\mathbf{a}_{*r}^{(N)}\|_1 \}$  need to be computed. For robustness, the weight  $w_k$  is replaced by  $|u_k| \|\mathbf{a}_{*r}^{(1)}\|_1 \|\mathbf{a}_{*r}^{(2)}\|_1 \dots \|\mathbf{a}_{*r}^{(N)}\|_1 / (L^{(1)} \times L^{(2)} \dots \times L^{(N)})$ . In each turn, two components is calculated, and divided by the corresponding length.

### 7.3 Compute the $x_{i_1, i_2, \dots, i_N}$

- In each round of [line 8 of Algorithm 3](#), only two indexes of  $\{i_1, i_2, \dots, i_N\}$  need to be sampled with particular distribution ( $\mathbf{p}_n = \mathbf{a}_{*r}^{(n)} / \|\mathbf{a}_{*r}^{(n)}\|_1$ ).
- Update values  $\Delta x$

$$\Delta x = \text{sgn}(a_{i_1 k}^{(1)} a_{i_2 k}^{(2)} \cdot a_{i_3 k}^{(3)} \cdots a_{i_N k}^{(N)} \cdot u_k \cdot u_{k'}) a_{i_1 k'}^{(1)} a_{i_2 k'}^{(2)} a_{i_3 k'}^{(3)} \cdots a_{i_N k'}^{(N)}$$

And only four of elements need to be calculated.

- Suppose  $k-1 = xND + yD + z$  and  $x+1 = y+1 = n$ , then only  $i_n$  need to be sampled with specific distribution. After sampled  $i_n, k'$ , we update  $x_{i_1, i_2, i_3, \dots, i_N}$  for  $i_m = \{1, \dots, L^{(m)}\}, m \neq n$

$$x_{i_1, i_2, i_3, \dots, i_N} \leftarrow x_{i_1, i_2, i_3, \dots, i_N} + \frac{1}{\prod_{m \neq n} L^{(m)}} \Delta x$$

- Or only  $i_n, i_m$  need to be sampled with specific distribution. After sampled  $i_n, i_m, k'$ , then we update  $x_{i_1, i_2, i_3, \dots, i_N}$  for  $i_h = \{1, \dots, L^{(h)}\}, h \neq m, n$

$$x_{i_1, i_2, i_3, \dots, i_N} \leftarrow x_{i_1, i_2, i_3, \dots, i_N} + \frac{1}{\prod_{h \neq m, n} L^{(h)}} \Delta x$$