

# Tensor link analysis

Bruno Fabulić, Helena Marciuš, Dora Parmać

6. ožujka 2021.

# Sadržaj

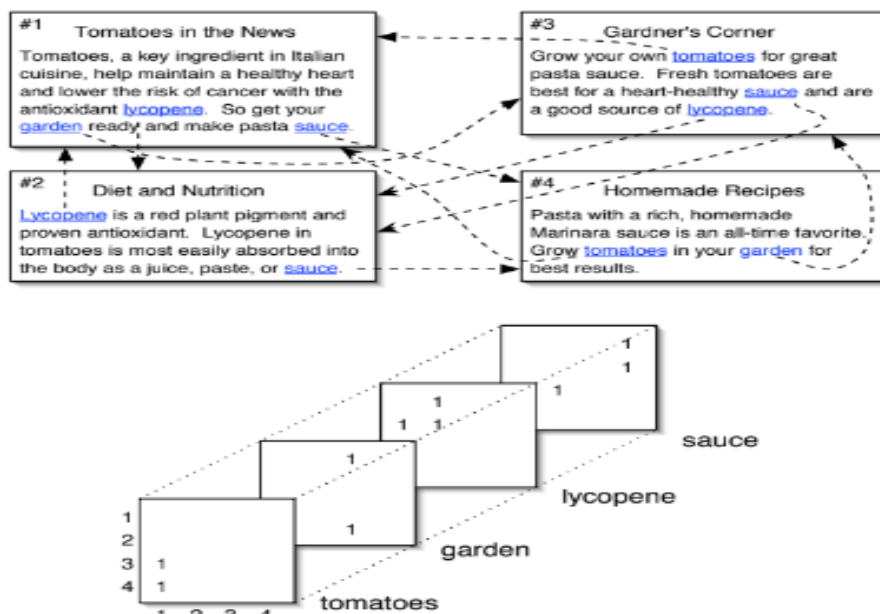
<b>1</b>	<b>Uvod</b>	<b>1</b>
<b>2</b>	<b>Usporedba HITS i TOPHITS algoritma</b>	<b>2</b>
2.1	HITS . . . . .	2
2.2	TOPHITS . . . . .	4
2.2.1	TOPHITS ukratko . . . . .	6
<b>3</b>	<b>Metodologija</b>	<b>7</b>
3.1	PARAFAC model . . . . .	7
3.2	Pohlepni PARAFAC . . . . .	7
3.2.1	Pohlepni PARAFAC - pseudokod . . . . .	8
3.3	Alternirajući najmanji kvadrati za PARAFAC . . . . .	9
3.3.1	Alternirajući najmanji kvadrati za PARAFAC - pse- udokod . . . . .	10
3.3.2	Inicijalizacija . . . . .	10
3.4	Specijalna konsideracija za razbacane podatke . . . . .	10
<b>4</b>	<b>TOPHITS i upiti</b>	<b>12</b>
4.1	Pronalaženje podudarne grupe . . . . .	12
4.2	Pronalaženje skupa autoriteta . . . . .	12
<b>5</b>	<b>Implementacija</b>	<b>14</b>

# 1 Uvod

Kako se veličina Weba povećava, postaje sve važnije analizirati strukturu linkova uzimajući u obzir i kontekst. Multilinearna algebra pruža novi alat za obuhvaćanje teksta poveznice i drugim informacija koji se koriste u drugih metodama analize poveznica kao što je HITS algoritam.

Algoritam kojeg proučavamo naziva se Topical HITS ili TOPHITS. On je generalizacija Kleinbergovog HITS modela. TOPHITS dodaje 3. dimenziju za oblikovanje tenzora susjedstva koji uključuje i tekst poveznice. Ova dodatna informacija pruža način za obuhvaćanje konteksta uz izračunavanje autoriteta i hubova, što je postignuto 3D Parallel Factors (PARAFAC) dekompozicijom, više-dimenzionalni analogon SVD dekompozicije. Dodavanje 3. dimenzije znatno poboljšava primjenjivost HITS-a, zato što se TOPHITS analiza može izvoditi unaprijed i offline. Također, TOPHITS otkriva latentno grupiranje stranice i informacija.

U ovom seminaru, opisat ćemo brži matematički algoritam za izračunavanje TOPHITS algoritma na podacima.



Slika 1: TOPHITS analizira 3D tenzor koji reprezentira kolekciju web stranica

## 2 Usporedba HITS i TOPHITS algoritma

Mnoge metode za analizu weba, poput PageRanka i HITS-a bazirane su na matrici susjedstva grafa kolekcije web stranica. PageRank rezultati dani su ulazima glavnog svojstvenog vektora Markovljeve matrice vjerojatnosti. S druge strane, HITS računa hubove i autoritete za svaki vrh i oni odgovaraju glavnom lijevom i desnom singularnom vektoru matrice susjedstva. Zanimljivo svojstvo HITS-a, što ne dijeli s PageRankom jest da se višestruki parovi singularnih vektora mogu razmatrati.

### 2.1 HITS

Neka je  $I$  broj stranica u našem podgrafu. Svaka stranica ima hub-score  $h$  i autoritet-score  $a$  koje računamo iterativno na sljedeći način:

$$\mathbf{h}_i^{(t+1)} = \sum_{i \rightarrow j} \mathbf{a}_j^t \quad (1)$$

$$\mathbf{a}_j^{(t+1)} = \sum_{i \rightarrow j} \mathbf{h}_i^{t+1} \quad (2)$$

Nakon svake iteracije, normaliziramo  $h$  i  $a$ . To znači da je hub-score stranice i jednak sumi autoritet-scorova svih stranica na koje pokazuje. Slično, autoritet-score stranice i jednak je sumi hub-scorova svih stranica koje pokazuju na nju.

Zapišimo to matrično. Neka je  $I \times I$  matrica susjedstva  $\mathbf{X}$  definirana kao:

$$X_{i,j} = \begin{cases} 1 & \text{ako stranica } i \text{ pokazuje na stranicu } j \\ 0 & \text{inače} \end{cases} \quad (3)$$

Tada jednačbe 1 i 2 postaju:

$$\mathbf{h}^{(t+1)} = \mathbf{X} \mathbf{a}^{(t)}, \mathbf{a}^{(t+1)} = \mathbf{X}^T \mathbf{h}^{(t+1)} \quad (4)$$

Prvih  $R$  faktora singularne dekompozicije of  $\mathbf{X}$  proizvodi najbolju  $R$ -rang aproksimaciju, odnosno možemo aproksimirati  $\mathbf{X}$  kao:

$$\mathbf{X} \approx \mathbf{H} \mathbf{\Sigma} \mathbf{A}^T \equiv \sum_{r=1}^R \sigma_r \mathbf{h}_r \circ \mathbf{a}_r \quad (5)$$

Vrijedi  $\mathbf{\Sigma} = \text{diag}\{\sigma_1 \dots \sigma_R\}$  i za singularne vrijednosti pretpostavljamo  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_R > 0$  i gdje su  $\mathbf{h}_r$  i  $\mathbf{a}_r$  odgovarajući singularni vektori. Matrice

$\mathbf{H}$  i  $\mathbf{A}$  su veličine  $I \times R$  i imaju ortonormalna polja. To možemo zamisliti kao aproksimaciju matrice  $\mathbf{X}$  sumom  $R$  1 ranga vanjskog produkta, kako je prikazano na (2) Glavni par singularnih vektora  $\mathbf{h}_1$  i  $\mathbf{a}_1$  pruža hub i aurotitet rezultate za dominantnu temu u kolekciji web stranica. Drugim riječima, stranice koje imaju najveće rezultate  $\mathbf{h}_1$  su najbolji hubovi za dominantnu temu, a stranice koje imaju najveće rezultate u  $\mathbf{a}_1$  su najbolji autoriteti.

$$\square = \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} + \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} \begin{array}{|c|} \hline \text{ } \\ \hline \end{array} + \dots$$

Slika 2: SVD dekompozicija

## 2.2 TOPHITS

Ilustrirajmo sada TOPHITS algoritam. On proizvede trojke  $\mathbf{h}_1, \mathbf{a}_1, \mathbf{t}_1$  gdje  $\mathbf{u}$  i  $\mathbf{v}$  predstavljaju hub i aurotitet rezultate, a  $\mathbf{t}$  term rezultate.

Slično kao HITS, rezultati se mogu dobiti iterativno. Neka je  $I$  broj stranica, a  $K$  broj termova. Tada je:

$$\mathbf{h}_i^{(t+1)} = \sum_{i \xrightarrow{k} j} \mathbf{a}_j^t \mathbf{t}_k^{(t)} \quad \text{za } i = 1, \dots, I \quad (6)$$

$$\mathbf{a}_j^{(t+1)} = \sum_{i \xrightarrow{k} j} \mathbf{h}_i^{t+1} \mathbf{t}_k^{(t)} \quad \text{za } j = 1, \dots, I \quad (7)$$

$$\mathbf{t}_k^{(t+1)} = \sum_{i \xrightarrow{k} j} \mathbf{a}_j^{t+1} \mathbf{h}_i^{t+1} \quad \text{za } k = 1, \dots, I \quad (8)$$

Kao i kod HITS-a, normaliziramo vrijednosti nakon svake iteracije.

Zapišimo to u tenzor formi. Neka je  $\mathbf{X}$   $I \times I \times K$  tenzor susjedstva podgrafa definiran kao:

$$X_{i,j,k} = \begin{cases} 1 & \text{ako stranica } i \text{ pokazuje na stranicu } j \text{ koristeći term } k \\ 0 & \text{inače} \end{cases} \quad (9)$$

Tada se jednadžbe 6, 7 i 8 mogu zapisati kao:

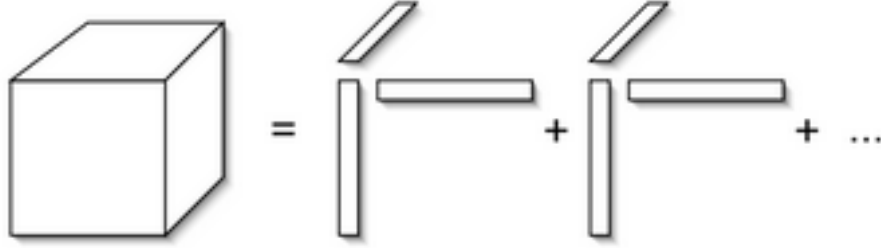
$$\mathbf{h}^{(t+1)} = \mathbf{X} \bar{\times}_2 \mathbf{a}^{(t)} \bar{\times}_3 \mathbf{t}^{(t)} \quad (10)$$

$$\mathbf{a}^{(t+1)} = \mathbf{X}^T \bar{\times} \mathbf{h}^{(t+1)} \bar{\times}_3 \mathbf{t}^{(t)} \quad (11)$$

$$\mathbf{t}^{(t+1)} = \mathbf{X}^T \bar{\times}_1 \mathbf{h}^{(t+1)} \bar{\times}_2 \mathbf{a}^{(t)} \quad (12)$$

Ovdje notacija  $\mathbf{X} \bar{\times}_i \mathbf{x}$  znači množenje tenzora  $\mathbf{X}$  vektorom  $\mathbf{x}$  u dimenziji  $i$ .

Tada se aproksimacija tenzora  $\mathbf{X}$  PARAFAC dekompozicijom može dobiti kao:



Slika 3: PARAFAC dekompozicija

$$\mathcal{X} \approx \lambda[[\mathbf{H}, \mathbf{A}, \mathbf{T}]] \equiv \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r \circ \mathbf{t}_r \quad (13)$$

Ovdje pretpostavljamo  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_R$ . Matrice  $\mathbf{H}, \mathbf{A}, \mathbf{T}$  imaju polja duljine 1, ali suprotno od SVD-a, nisu generalno ortonormalne. PARAFAC dekompozicija aproksimira tenzor  $\mathcal{X}$  sumom, kao što je prikazano na (3).

Glavna trojka PARAFAC vektora  $\mathbf{h}_1, \mathbf{a}_1, \mathbf{t}_1$  pruža hub i aurotitet rezultate za dominantnu temu u kolekciji web stranica (na isti način kao što je prethodno opisano kod HITS-a), te term rezultate. Drugim riječima, termi koji imaju najveće rezultate  $\mathbf{t}_1$  su najdeskriptivniji termi.

Suprotno SVD-u, nema garancije da će PARAFAC aproksimacija biti optimalna, te vektori nisu ortogonalni.

### 2.2.1 TOPHITS ukratko

Ideja iza TOPHITSa je sljedeća: pretpostavimo da analiziramo kolekciju  $I$  web stranica imajući ukupno  $K$  termova u tekstu poveznica svih poveznica. Tada je  $I \times I \times K$  tenzor susjedstva  $\mathcal{X}$  definiran kao u (6). Primijetimo da je tenzor  $\mathcal{X}$  rijedak zato što većina stranica pokazuje samo na nekoliko drugih stranica u kolekciji, a svaki link koristi samo nekoliko termova. Stoga je realno za očekivati da je broj ne-nul elemenata u  $\mathcal{X}$  jednak  $O(I)$ .

Uzimajući u obzir  $R > 0$  (broj različitih grupiranja podataka), TOPHITS algoritam pronađe matrice  $\mathbf{H} \mathbf{i} \mathbf{A}$ , veličina  $I \times R$ , i matricu  $\mathbf{T}$  veličine  $K \times R$ . Svaka trojka  $\{\mathbf{h}_r, \mathbf{a}_r, \mathbf{t}_r\}$ , za  $r = 1 \dots R$  definira grupiranje hubova, autoriteta i terma uzimajući u obzir ulaze s najvećom vrijednosti u svakom vektoru; vrijednost  $\lambda_r$  definira težinu grupiranja. (Bez smanjenja općenitosti, pretpostavljamo da su stupci naših matrica normalizirani.)

TOPHITS je zapravo ekstenzija Kleinbergovog HITS algortma koja je koristila singularni vektor matrice hiperlinkova (2D tenzor) da bi proizvela skup hubova i autoriteta. Dodatak topic-vektora znači da je određivanje skupa singularnih vektora koji sadrže odgovor na upit zapravo promatranje koji topic-vektori imaju najveći score uključujući njihove autoritete i hubove. Slično kao PageRank, TOPHITS je neovisan o upitima zato što se računanje značajnih vektora može obaviti unaprijed i offline.



### 3 Metodologija

#### 3.1 PARAFAC model

3D dekompozicija simultano je predložena sa strane Harshmana, koristeći ime "Parallel Factors ili PARAFAC" te sa strane Carroll i Chang, koristeći ime "Canonical Decomposition ili CANDECOMP". Cilj je prikazati vektor reda  $N$  kao sumu vektora vanjskog produkta kao što je prikazano na slici 3.

Matematički, problem je zadan na sljedeći način. Pretpostavimo da imamo tenzor  $\mathcal{X}$  veličine  $I_1 \times I_2 \times \dots \times I_n$  i željeni rang aproksimacije  $R$ . Tada želimo naći matrice  $U^n$  veličine  $I_n \times R$ , za svaki  $n = 1, \dots, N$  i težinski vektor  $\lambda$  duljine  $R$  tako da vrijedi:

$$\mathcal{X} \approx \lambda[[U^{(1)}, U^{(2)}, \dots, U^{(N)}]] \quad (14)$$

Kruskalov operator  $[[\cdot]]$  je skraćeni zapis sume ranga svih vanjskih produkata stupaca. Drugim riječima,

$$\lambda[[U^{(1)}, U^{(2)}, \dots, U^{(N)}]] \equiv \sum_{r=1}^R \lambda_r u_r^{(1)} \circ u_r^{(2)} \circ \dots \circ u_r^{(N)} \quad (15)$$

Bez smanjenja općenitosti, prepostavljamo da je  $\|u_r^{(n)}\| = 1$  za sve  $r = 1, \dots, R$  i  $n = 1, \dots, N$ . Povrh toga, preuredimo završno rješenje tako da vrijedi  $\lambda_1 \geq \lambda_2 \geq \dots \lambda_R$ .

Cilj nam je riješiti problem minimizacije:

$$\min \left\| \mathcal{X} - \lambda[[U^{(1)}, U^{(2)}, \dots, U^{(N)}]] \right\|^2 \quad (16)$$

s obzirom na  $\lambda \in \mathbb{R}^R$ ,  $U^{(n)} \in \mathbb{R}^{I_n \times R}$ ,  $n = 1, \dots, N$ .

U slučaju TOPHITSA,  $\lambda$  je tenzor reda 3, tj  $N = 3$  i vrijedi:

$$H \equiv U^{(1)}, A \equiv U^{(2)}, T \equiv U^{(3)} \quad (17)$$

#### 3.2 Pohlepni PARAFAC

Notacija  $\mathbf{X}_{(n)}$  reprezentira  $n$ -to unfolding tenzora  $\mathcal{X}$ . Drugim riječima,  $\mathbf{X}_{(n)}$  predstavlja reorganizirane ulaze  $\mathcal{X}$  u matricu veličine  $I_n \times J$ , gdje je  $J = \prod_{k=1, k \neq n}^N I_k$  tako da su "vlakna" u dimenziji  $n$  uređena kao stupci matrice. Matematički imamo:

$$\begin{aligned}
[\mathbf{X}_{(n)}]_{i,j} &= x_{i_1, i_2, \dots, i_N} \\
i &= i_n, j = 1 + \sum_{k=1, k \neq n}^N (i_n - 1) \prod_{l=1, l \neq n^{k-1}} I_l \\
i &\leq i \leq I_n, 1 \leq j \leq J \quad (18)
\end{aligned}$$

### 3.2.1 Pohlepni PARAFAC - pseudokod

---

**Algorithm 1** Greedy PARAFAC

---

**in:** Tensor  $\mathfrak{X}$  of size  $I_1 \times I_2 \times \dots \times I_N$ .  
**in:** Desired rank  $R > 0$ .  
**for**  $r = 1, \dots, R$  **do** {outer loop}  
    Set  $\mathbf{v}^{(n)}$  to be a vector of all ones of length  $I_n$  for  $n = 1, \dots, N$ .  
    **repeat** {middle loop}  
        **for**  $n = 1, \dots, N$  **do** {inner loop}  
            Set  $\mathbf{w} = \mathbf{X}_{(n)} \mathbf{z}^{(n)} - \sum_{i=1}^{r-1} \left( \mathbf{u}_i^{(n)} \prod_{\substack{m=1 \\ m \neq n}}^N (\mathbf{v}^{(m)})^\top \mathbf{u}_i^{(m)} \right)$  where  $\mathbf{z}^{(n)} \equiv \mathbf{v}^{(1)} \otimes \dots \otimes \mathbf{v}^{(n-1)} \otimes \mathbf{v}^{(n+1)} \otimes \dots \otimes \mathbf{v}^{(N)}$ .  
            Set  $\lambda_r = \|\mathbf{w}\|$ .  
            Set  $\mathbf{v}^{(n)} = \mathbf{w} / \lambda_r$ .  
        **end for**  
    **until** the fit ceases to improve or the maximum number of middle-loop iterations has been exceeded.  
    Set  $\mathbf{u}_r^{(n)} = \mathbf{v}^{(n)}$  for  $n = 1, \dots, N$ .  
**end for**  
**out:**  $\lambda \in \mathbb{R}^R$  and  $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R}$  for  $n = 1, \dots, N$ .

---

1. Vanjska petlja (po  $r$ ) računa  $\{u_r^{(1)}, \dots, u_r^{(N)}\}$   
Za svaki  $r, r = 1, \dots, N$ , unutarnja petlja metodom alternirajućih najmanjih kvadrata minimizira

2. 
$$\left\| \left( X - \sum_{i=1}^{r-1} \lambda_i u_i^{(1)} \circ \dots \circ u_i^{(N)} \right) - \left( v^{(1)} \circ \dots \circ v^{(N)} \right) \right\| \quad (19)$$

po vektorima  $v^{(1)}, \dots, v^{(N)}$

Tada je  $u_r^{(n)} = v^{(n)}$ ,  $n = 1, \dots, N$

### 3.3 Alternirajući najmanji kvadrati za PARAFAC

Češći pristup rješavanja PARAFAC dekompozicija je metoda alternirajućih najmanjih kvadrata. U svakoj unutarnjoj iteraciji, računamo cijelu ntu matricu  $\mathbf{U}^{(n)}$  dok sve ostale matrice držimo fiksne.  $V$  koji se računa nakon svake unutarnje iteracije je rješenje danog optimizacijskog problema:

$$\min_V \left\| X - [[U^{(1)}, \dots, U^{(n-1)}, V, U^{(n+1)}, \dots, U^{(N)}]] \right\|^2 \quad (20)$$

što se matrično može zapisati kao

$$\min_V \left\| X_n - VZ^{(n)T} \right\|^2 \quad (21)$$

Matrica  $\mathbf{Z}^{(n)}$  je veličine  $J \times R$  i definirana kao (2.6) Rješenje najmanjih kvadrata za 21 uključuje pseudo inverz matrice  $\mathbf{Z}^{(n)}$ :

$$\mathbf{V} = \mathbf{X}_{(n)} (\mathbf{Z}^{(n)T})^\dagger \quad (22)$$

Prikladno, pseudo inverz  $\mathbf{Z}^{(n)}$  ima specijalnu strukturu. Neka je  $\mathbf{Y}^{(n)}$  simetrična  $R \times R$  matrica kao u (2.7). Tada se može pokazati da je:

$$(\mathbf{Z}^{(n)T})^\dagger = \mathbf{Z}^{(n)} \mathbf{Y}^{(n)T} \quad (23)$$

Dakle rješenje za 21 je dano (2.5), pa računanje  $\mathbf{U}_n$  zapravo svodimo na računanje inverza specijalne  $R \times R$  matrice  $\mathbf{Y}^{(n)}$ .

### 3.3.1 Alternirajući najmanji kvadrati za PARAFAC - pseudokod

---

**Algorithm 2** Alternating Least Squares (ALS) for N-way arrays

---

**in:** Tensor  $\mathcal{X}$  of size  $I_1 \times I_2 \times \dots \times I_N$ .  
**in:** Desired rank  $R > 0$ .  
Initialize  $\mathbf{U}^{(n)}$  for  $n = 1, \dots, N$  (see §2.4).  
**repeat** {outer loop}  
    **for**  $n = 1, \dots, N$  **do** {inner loop}

$$(2.5) \quad \text{Set } \mathbf{V} = \mathbf{X}_{(n)} \mathbf{Z}^{(n)} \mathbf{Y}^{(n)},$$

$$(2.6) \quad \text{where } \mathbf{Z}^{(n)} \equiv \sum_{r=1}^R \mathbf{u}_r^{(1)} \otimes \dots \otimes \mathbf{u}_r^{(n-1)} \otimes \mathbf{u}_r^{(n+1)} \otimes \dots \otimes \mathbf{u}_r^{(N)},$$

$$(2.7) \quad \text{and } \mathbf{Y}^{(n)} \equiv \left( \mathbf{U}^{(1)\top} \mathbf{U}^{(1)} * \dots * \mathbf{U}^{(n-1)\top} \mathbf{U}^{(n-1)} * \mathbf{U}^{(n+1)\top} \mathbf{U}^{(n+1)} * \dots * \mathbf{U}^{(N)\top} \mathbf{U}^{(N)} \right)^{-1}.$$

**for**  $r=1, \dots, R$  **do** {Assign  $\mathbf{U}^{(n)}$ }

        Set  $\lambda_r = \|\mathbf{v}_r\|$

        Set  $\mathbf{u}_r^{(n)} = \mathbf{v}_r / \lambda_r$ .

**end for**

**until** the fit ceases to improve or the maximum number of outer iterations is exceeded.

**out:**  $\lambda \in \mathbb{R}^R$  and  $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times R}$  for  $n = 1, \dots, N$ .

---

### 3.3.2 Inicijalizacija

- Pohlepni PARAFAC - inicijaliziramo matrice  $\mathbf{U}^{(n)}$  kao matrice dobivene pohlepnim PARAFAC algoritmom
- Nasumična inicijalizacija - matrice  $\mathbf{U}^{(n)}$  inicijaliziramo nasumičnim vrijednostima
- HOSVD inicijalizacija - za svaki mod  $\mathbf{X}_{(n)}$  izračunamo SVD dekompoziciju,  $\mathbf{X}_{(n)} = \mathbf{U}_n \Sigma_n \mathbf{V}_n^T$ , i inicijaliziramo matrice  $\mathbf{U}^{(n)} = \mathbf{U}_n$ ,  $n = 1, \dots, N$

### 3.4 Specijalna konsideracija za razbacane podatke

Kao što smo već spomenuli, tenzor  $\mathcal{X}$  je veoma rijedak. Posljedično, njegova reprezentacija  $\mathbf{X}_n$  je rijetka matrica (ima iste ne nul elemente, ali preoblikovane). Matrica  $\mathbf{Z}^{(n)}$  iz (2.6) se ne bi trebala formirati eksplicitno jer bi onda bila gusta matrice veličine  $I_n \times J$  gdje je  $J = \prod_{k=1, k \neq n}^N I_k$ . Umjesto toga, računanje

$$\mathbf{X}_{(n)} \mathbf{Z}^{(n)} \tag{24}$$

potrebno za (2.5) mora biti izračunato posebno, iskorištavanjem svojstava Kroneckerove strukture produkta u  $\mathbf{Z}^{(n)}$  da bi se očuvala rijetkost. Završni rezultat je veličine  $I_n \times R$  stoga sse može pohraniti u gustu matricu. Jedna od metoda prikazana je sljedećim algoritmom:

---

**Algorithm 3** Computing the sparse product  $\mathbf{X}_{(n)}\mathbf{Z}^{(n)}$

---

**in:** Tensor  $\mathbf{X}$  of size  $I_1 \times I_2 \times \dots \times I_N$  with  $Q$  nonzeros.  
Let the index of the  $q$ th nonzero be  $(k_{1_q}, k_{2_q}, \dots, k_{N_q})$   
and its value be given by  $v_q$ .  
**in:** Index  $n$  and matrices  $\mathbf{U}^{(m)}$  for  $1 \leq m \leq N, m \neq n$ .  
**for**  $r = 1 \dots, R$  **do**  
    **for**  $q = 1, \dots, Q$  **do**  
        Compute  $w_q = v_q \prod_{\substack{m=1 \\ m \neq n}}^N u_{k_{m_q}, r}^{(m)}$   
    **end for**  
    **for**  $i = 1, \dots, I_n$  **do** {Compute  $r$ th column of  $\mathbf{P}$ }  
        Set  $p_{ir} = \sum_{\substack{q=1 \\ k_{n_q} = i}}^Q w_q$ .  
    **end for**  
**end for**  
**out:**  $\mathbf{P} = \mathbf{X}_{(n)}\mathbf{Z}^{(n)}$

---

## 4 TOPHITS i upiti

Jednom kada smo izračunali TOPHITS model ranga  $R$

$$\mathcal{X} = \lambda[H, A, T] \quad (25)$$

možemo ga iskoristiti za razumijevanje podataka na razne načine. Gledajući najveću vrijednost svake trojke  $\{\mathbf{h}_r, \mathbf{a}_r, \mathbf{t}_r\}$  pruža grupiranje hubova, autoriteta, termova web stranice, a  $\lambda_r$  pruža relativnu težinu grupiranja.

Jedno pitanje koje smatramo osnovnim pitanjem web pretraživanja: pronaći sve stranice povezane s određenim termom ili skupom termova. Promotrimo vektor upita  $g$  duljine  $K$  (gdje je  $K$  broj termova) kao:

$$q_k = \begin{cases} 1 & \text{ako je term } k \text{ u upitu} \\ 0 & \text{inače} \end{cases} \quad (26)$$

Primijetimo da nema razloga restringirati se na upite o termovima. Također se možemo baviti problem pronalaženja web stranica i/ili termova povezanih s određenom stranicom ili skupom stranica.

### 4.1 Pronalaženje podudarne grupe

Umjesto vraćanja liste rangiranih stranica, TOPHITS omogućuje opciju identificiranja grupiranja relevantnih za dani upit. Možemo kreirati vektor grupe  $s$  duljine  $R$  koji sadrži vrijednosti svakog grupiranja, na temelju matrice  $T$  iz PARAFAC modela:

$$s = \Lambda T^T q, \Lambda = \text{diag}(\lambda) \quad (27)$$

Ulaz  $s_r$  daje rezultat-vrijednost  $r$ -te grupe, a grupiranja s većim rezultatima smatraju se više relevantnima.

Alternativno, možemo konstruirati vektora upita na temelju web stranica,  $\hat{q} \in R^l$  te izračunati rezultate grupe kao:

$$\hat{s} = \Lambda A^T \hat{q}, \Lambda = \text{diag}(\lambda) \quad (28)$$

### 4.2 Pronalaženje skupa autoriteta

Također je moguće vratiti tradicionalnu rangiranu listu vjerojatnosti. Možemo kombinirati sve informacije TOPHITS modela da bi vratili skup rangiranih autoriteta i/ili hubova. Definirajmo  $s$  kao u 27. Tada su kombinirani autoriteti dani s:

$$a^* = As = \sum_{r=1}^R s_r a_r \quad (29)$$

Sortiranje ulaza  $\mathbf{a}^*$  vraća rangiranu listu autoriteta. Na isti način, kombinirani hubovi dani su:

$$h^* = Hs = \sum_{r=1}^R s_r h_r \quad (30)$$

## 5 Implementacija

Testirali smo našu tehniku na skupu podataka, generiranim koristeći web crawler koji uključuje i tekst poveznice.