

# TOPHITS algoritam

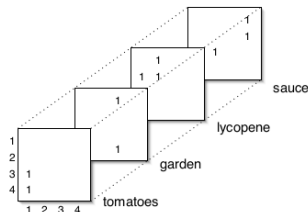
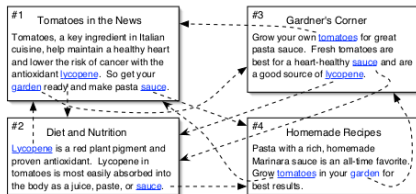
Bruno Fabulić, Helena Marcioš, Dora Parmač

8. ožujka 2021.

# Sadržaj

- 1 HITS algoritam
- 2 TOPHITS algoritam
- 3 PARAFAC dekompozicija
  - Pohlepni PARAFAC
  - Alternirajući najmanji kvadrati za PARAFAC
- 4 TOPHITS i upiti

- Povećanjem *Weba*, sve je važnije analizirati strukturu linkova uzimajući u obzir i kontekst
- TOPHITS algoritam - generalizacija HITS algoritma, dodaje treću dimenziju za oblikovanje tenzora susjedstva koji uključuje tekst poveznice



# HITS algoritam

- Internet prikazujemo usmjerenim grafom  $G$  - *web* stranice su čvorovi, poveznice između stranica su bridovi
- Postoji brid između čvora  $i$  i čvora  $j$  ako postoji poveznica sa stranice  $i$  na stranicu  $j$
- Svakoj *web* stranici pridružuje se autoritet-vrijednost  $a_i$  te hub-vrijednost  $h_i$
- Neka je  $n$  broj *web* stranica. Autoritet-vrijednost  $i$  hub-vrijednost stranice  $i$  definira se sljedećim formulama:

$$a_i = \sum_{j \rightarrow i} h_j \quad h_i = \sum_{i \rightarrow j} a_j, \quad i = 1, \dots, n \quad (1)$$

- Ove se vrijednosti iterativno računaju na sljedeći način:

$$h_i^{(t+1)} = \sum_{i \rightarrow j} a_j^{(t)} \quad a_i^{(t+1)} = \sum_{j \rightarrow i} h_j^{(t+1)}, \quad i = 1, \dots, n \quad (2)$$

# HITS algoritam

- Usmjerenom grafu  $G$  pridružujemo matricu susjedstva  $A$ :

$$A_{ij} = \begin{cases} 1, & \text{ako } i \rightarrow j \\ 0, & \text{inače} \end{cases} \quad (3)$$

- Jednakosti (2) mogu se zapisati kao

$$h^{(t+1)} = Aa^{(t)} \quad a^{(t+1)} = Ah^{(t+1)} \quad (4)$$

gdje su  $h$  i  $a$  vektori hub, tj. autoritet vrijednosti

- Koristeći SVD dekompoziciju, možemo pronaći aproksimaciju od  $A$  ranga  $p$ :

$$A \approx \sum_{i=1}^p \sigma^{(i)} u^{(i)} \circ v^{(i)} \quad (5)$$

- Vektori hub i autoritet vrijednosti u iteracijama (4) konvergiraju dominantnim singularnim vektorima:

$$h^{(t)} \rightarrow u^{(1)} \quad a^{(t)} \rightarrow v^{(1)} \quad (6)$$

# TOPHITS algoritam

- Kombinira *anchor* text s *hyperlink* strukturom Interneta
- Osim hub i autoritet vrijednosti, računa i *topic* vrijednost za pojmove u *anchor* tekstu
- Može se računati iterativno. Ako je  $n$  broj *web* stranica i  $m$  broj pojmova. Tada vrijedi

$$\begin{aligned}h_i^{(t+1)} &= \sum_{i \xrightarrow{k} j} a_j^{(t)} t_k^{(t)} \\a_j^{(t+1)} &= \sum_{i \xrightarrow{k} j} h_i^{(t+1)} t_k^{(t)} \\t_k^{(t+1)} &= \sum_{i \xrightarrow{k} j} a_j^{(t+1)} h_i^{(t+1)}\end{aligned}\tag{7}$$

# TOPHITS algoritam

- Za prikazivanje veza između *web* stranica koristi se tenzor  $\mathbf{A}$  reda 3:

$$A_{ijk} = \begin{cases} 1, & \text{ako } i \rightarrow j \text{ s anchor tekstem } k \\ 0, & \text{inače} \end{cases} \quad (8)$$

- Jednakosti (7) mogu se zapisati kao

$$\begin{aligned} h^{(t+1)} &= \mathbf{A} \times_2 a^{(t)} \times_3 t^{(t)} \\ a^{(t+1)} &= \mathbf{A} \times_1 h^{(t+1)} \times_3 t^{(t)} \\ t^{(t+1)} &= \mathbf{A} \times_1 h^{(t+1)} \times_2 a^{(t+1)} \end{aligned} \quad (9)$$

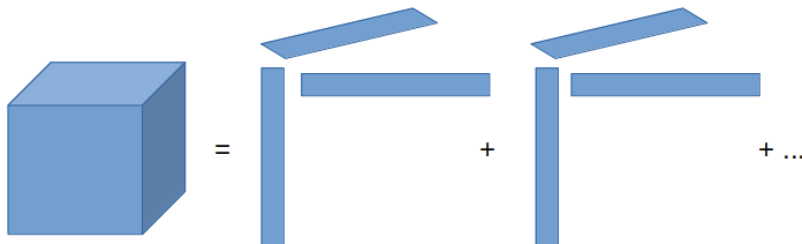
- $\times_n$  - množenje u modu  $n$ . Npr., za  $h = \mathbf{A} \times_2 a \times_3 t$  vrijedi

$$h_i = \sum_{j=1}^n \sum_{k=1}^m A_{ijk} a_j t_k, \quad i = 1, \dots, n$$

# TOPHITS algoritam

- Za izračunavanje vektora hub, autoritet i topic vrijednosti, može se koristiti PARAFAC dekompozicija
- Pomoću PARAFAC dekompozicije, računamo aproksimaciju ranga  $p$  tenzora  $\mathbf{A}$ :

$$\mathbf{A} \approx \sum_{i=1}^p \sigma^{(i)} \mathbf{u}^{(i)} \circ \mathbf{v}^{(i)} \circ \mathbf{w}^{(i)} \quad (10)$$





# TOPHITS algoritam

- Može se implementirati algoritam koji uz odgovarajuće uvjete daje aproksimaciju tako da iteracije (9) konvergiraju prema dominantnim singularnim vektorima:

$$h^{(t)} \rightarrow u^{(1)}, \quad a^{(t)} \rightarrow v^{(1)}, \quad t^{(t)} \rightarrow w^{(1)} \quad (11)$$

- Najveća vrijednost u  $w^{(1)}$  definira dominantnu temu, a najveće vrijednosti u  $u^{(1)}$  i  $v^{(1)}$  dominantan hub i autoritet za dominantnu temu.

# PARAFAC dekompozicija

- Cilj: za tenzora reda  $N$  pronaći dekompoziciju tako da se tenzor može zapisati kao suma produkta vektora
- Neka je  $X$  tenzor reda  $N$  dimenzija  $I_1 \times I_2 \times \dots \times I_N$  i  $R > 0$  željeni rang
- Tražimo matrice  $U^{(n)}$  dimenzija  $I_n \times R$  za  $n = 1, \dots, N$  i vektor težina  $\lambda$  duljine  $R$  tako da

$$X \approx \lambda[[U^{(1)}, U^{(2)}, \dots, U^{(N)}]] = \sum_{r=1}^R \lambda_r u_r^{(1)} \circ u_r^{(2)} \circ \dots \circ u_r^{(N)} \quad (12)$$

- Rješavamo sljedeći problem minimizacije:

$$\min \left\| X - \lambda[[U^{(1)}, U^{(2)}, \dots, U^{(N)}]] \right\|^2 \quad (13)$$

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

# Pohlepni PARAFAC - pseudokod

**ulaz:** tenzor  $\mathbf{X}$  dimenzija  $I_1 \times I_2 \times \dots \times I_N$ , rang  $R > 0$

**for**  $r = 1, \dots, R$  **do**

inicijaliziraj  $\mathbf{v}^{(n)}$  kao vektor jedinica duljine  $I_n$ ,  $n = 1, \dots, N$

**do**

**for**  $n = 1, \dots, N$  **do**

$$\mathbf{z}^{(n)} = \mathbf{v}^{(1)} \otimes \dots \otimes \mathbf{v}^{(n-1)} \otimes \mathbf{v}^{(n+1)} \otimes \mathbf{v}^{(N)}$$

$$\mathbf{w} = \mathbf{X}_{(n)} \mathbf{z}^{(n)} - \sum_{i=1}^{r-1} \left( \mathbf{u}_i^{(n)} \prod_{\substack{m=1 \\ m \neq n}}^N (\mathbf{v}^{(m)})^T \mathbf{u}_i^{(m)} \right)$$

$$\lambda_r = \|\mathbf{w}\|$$

$$\mathbf{v}^{(n)} = \mathbf{w} / \lambda_r$$

**end for**

**while** fit ceases to improve

$$\mathbf{u}_r^{(n)} = \mathbf{v}^{(n)}, \quad n = 1, \dots, N$$

**end for**

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

# Pohlepni PARAFAC

- 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

$$\left\| \left( \mathbf{X} - \sum_{i=1}^{r-1} \lambda_i u_i^{(1)} \circ \dots \circ u_i^{(N)} \right) - \left( v^{(1)} \circ \dots \circ v^{(N)} \right) \right\| \quad (14)$$

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

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

# Alternirajući najmanji kvadrati za PARAFAC - pseudokod

**ulaz:** tenzor  $\mathbf{X}$  dimenzija  $I_1 \times I_2 \times \dots \times I_N$ , rang  $R > 0$

Inicijaliziraj  $U^{(n)}$  za  $n = 1, \dots, N$

**do**

**for**  $n = 1, \dots, N$  **do**

$$Z^{(n)} = U^{(1)} \odot \dots \odot U^{(n-1)} \odot U^{(n+1)} \odot \dots \odot U^{(N)}$$

$$Y^{(n)} = (U^{(1)T} U^{(1)} * \dots * U^{(n-1)T} U^{(n-1)} * \\ U^{(n+1)T} U^{(n+1)} * \dots * U^{(N)T} U^{(N)})^{-1}$$

$$V = X_{(n)} Z^{(n)} Y^{(n)}$$

**for**  $r = 1, \dots, R$  **do**

$$\lambda_r = \|v_r\|$$

$$u_r^{(n)} = v_r / \lambda_r$$

**end for**

**end for**

**while** fit ceases to improve

**izlaz:**  $\lambda \in \mathbb{R}^R$ ,  $U^{(n)} \in \mathbb{R}^{I_n \times R}$ ,  $n = 1, \dots, N$

# Alternirajući najmanji kvadrati za PARAFAC

- U svakoj iteraciji računamo matricu  $U^{(n)} \equiv V$ , a ostale su fiksirane, tj. rješavamo minimizacijski problem:

$$\min_V \left\| \mathbf{X} - [[U^{(1)}, \dots, U^{(n-1)}, V, U^{(n+1)}, \dots, U^{(N)}]] \right\|^2 \quad (15)$$

što se može zapisati kao

$$\min_V \left\| \mathbf{X}_n - \mathbf{V} \mathbf{Z}^{(n)T} \right\|^2 \quad (16)$$

- Rješenje gornjeg problema je

$$\mathbf{V} = \mathbf{X}_{(n)} \left( \mathbf{Z}^{(n)T} \right)^\dagger \quad (17)$$

- Može se pokazati

$$\left( \mathbf{Z}^{(n)T} \right)^\dagger = \mathbf{Z}^{(n)} \mathbf{Y}^{(n)T} \quad (18)$$

# Inicijalizacija $U^{(n)}$

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

## TOPHITS + PARAFAC

- Tenzor  $\mathbf{A}$  reda 3 dimenzija  $N \times N \times M$  i  $R > 0$
- PARAFAC dekompozicijom dobivamo matrice  $U^{(1)} \equiv H$  i  $U^{(2)} \equiv A$  dimenzija  $N \times R$  te  $U^{(3)} \equiv T$  dimenzije  $M \times R$ , te vektor  $\lambda \in \mathbb{R}^R$  tako da vrijedi

$$\mathbf{A} \approx \lambda[[H, A, T]] = \sum_{r=1}^R \lambda_r h_r \circ a_r \circ t_r \quad (19)$$

- Svaka trojka  $\{h_r, a_r, t_r\}$  određuje grupu hubova, autoriteta i pojmova;  $\lambda_r$  definira težinu grupe
- $\{h_1, a_1, t_1\}$  određuju najbolji hub, autoritet i pojam za dominantnu temu



# TOPHITS i upiti

- Kako pronaći sve stranice povezane termom ili skupom termova?
- Neka je  $q$  vektor duljine  $M$  ( $M$  je broj termova) definiran s

$$q_k = \begin{cases} 1, & \text{ako je term } k \text{ u upitu} \\ 0, & \text{inače} \end{cases} \quad (20)$$

- Pomoću ovog vektora možemo pronaći podudarne grupe i skupove autoriteta za zadani upit

# Pronalaženje podudarne grupe

- Možemo identificirati grupiranja relevantna za dani upit
- Računamo vektor  $s$  koji sadrži vrijednost svakog grupiranja:

$$s = \Lambda T^T q, \quad \Lambda = \text{diag}(\lambda) \quad (21)$$

- $s_r$  daje vrijednost  $r$ -te grupe, a što je vrijednost veća, grupa je relevantnija za dani upit
- Alternativno, možemo konstruirati vektor upita  $\hat{q} \in \mathbb{R}^N$  na temelju web stranica te izračunati rezultate grupe kao

$$\hat{s} = \Lambda A^T \hat{q}, \quad \Lambda = \text{diag}(\lambda) \quad (22)$$

# Pronalaženje skupa autoriteta

- Možemo vratiti tradicionalnu ranginaru listu vjerojatnosti - kombiniramo informacije dobivene TOPHITS-om da bi vratili skup kombiniranih autoriteta ili hubova
- Ako definiramo  $s$  kao u (21), tada su kombinirani autoriteti definirani  $s$

$$a^* = As = \sum_{r=1}^R s_r a_r \quad (23)$$

- Sortiranjem  $a^*$  dobivamo rangiranu listu autoriteta
- Analogno računamo kombinirane hubove:

$$h^* = Hs = \sum_{r=1}^R s_r h_r \quad (24)$$