

Tensor link analysis

Bruno Fabulić, Helena Marciuš, Dora Parmać

5. ožujka 2021.

Sadržaj

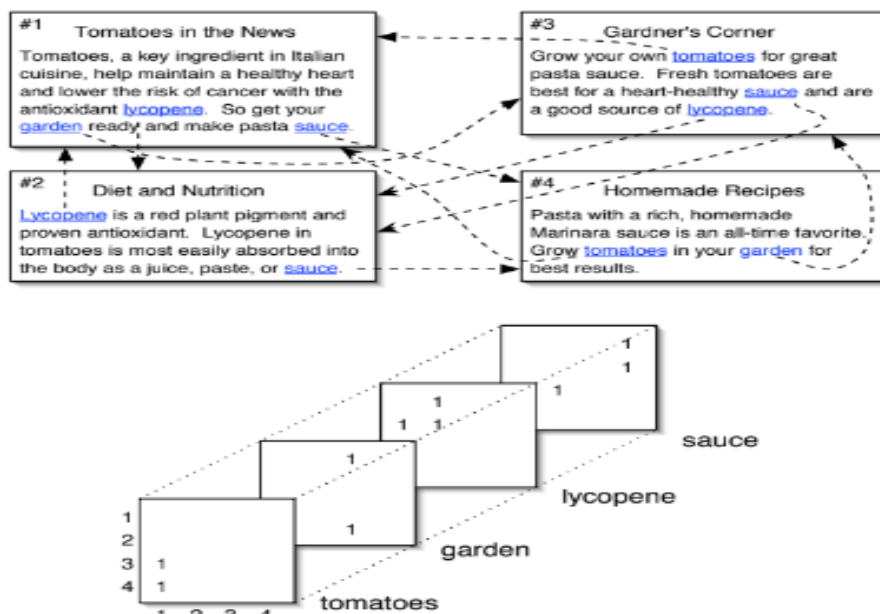
1	Uvod	1
2	HITS i TOPHITS	2
3	TOPHITS model	4
3.1	PARAFAC model	4
3.2	Pohlepni PARAFAC	5
3.2.1	Pohlepni PARAFAC - pseudokod	5
3.3	Alternirajući najmanji kvadrati za PARAFAC	6
3.3.1	Alternirajući najmanji kvadrati za PARAFAC - pse- udokod	8
3.3.2	Inicijalizacija	8
3.4	Specijalna konsideracija za razbacane podatke	8
4	TOPHITS i upiti	9
4.1	Pronalaženje podudarne grupe	10
4.2	Pronalaženje skupa autoriteta	10
5	Implementacija	10

1 Uvod

Kako se veličina Weba povećava, postaje sve važnije analizirati strukturu linka uzimajući u obzir i kontekst. Multilinearna algebra pruža novi alat za obuhvaćanje teksta poveznice i drugih 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 HITS i TOPHITS

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.

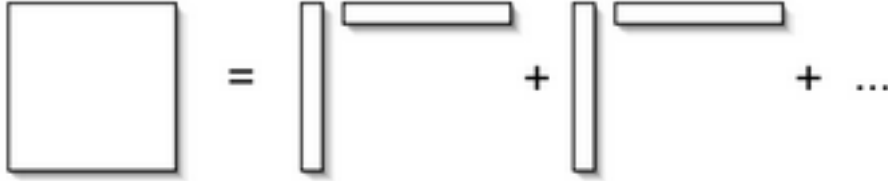
Promotrimo kolekciju I web stranica. U HITS-u, $I \times I$ matrica susjedstva \mathbf{X} definirana je kao:

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

HITS metoda može biti objašnjena na sljedeći način: koristi SVD matricnu dekompoziciju za računanje ranga R aproksimacije od \mathbf{X} :

$$\mathbf{X} \approx \mathbf{H}\Sigma\mathbf{A}^T \equiv \sum_{r=1}^R \sigma_r \mathbf{h}_r \circ \mathbf{a}_r \quad (2)$$

Vrijedi $\Sigma = \text{diag}\{\sigma_1 \dots \sigma_R\}$ i pretpostavljamo $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_R > 0$. 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 RANK-1 OUTER PRODUCTS, kako je prikazano na (2). Glavni par singularnih vektora \mathbf{h}_1 i \mathbf{a}_1 pruža hub i autoritet 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.



Slika 2: SVD dekompozicija

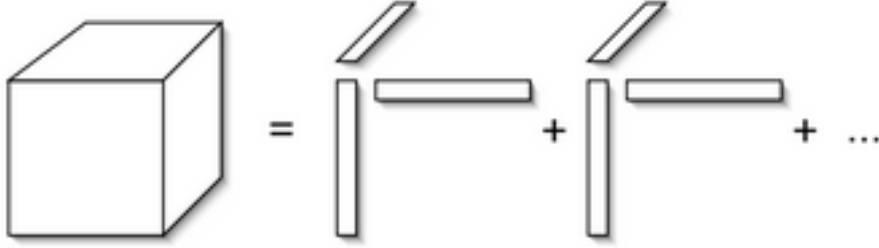
Ilustrirajmo sada TOPHITS algoritam. Tu je $I \times I \times K$ tenzor susjedstva \mathcal{X} definiran kao:

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

TOPHITS koristi PARAFAC model za generiranje R-ranga aproksimacije:

$$\mathcal{X} \approx \lambda[\mathbf{H}, \mathbf{A}, \mathbf{T}] \equiv \sum_{r=1}^R \lambda_r \mathbf{h}_r \circ \mathbf{a}_r \circ \mathbf{t}_r \quad (4)$$

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 R ranga 1 OUTER PRODUCTS, kao što je prikazano na (3)



Slika 3: PARAFAC dekompozicija

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.

3 TOPHITS model

Ideja iza TOPHITSa je sljedeća: pretpostavimo da analiziramo kolekciju I web stranica imajući ukupno K termova u tekstu poveznica svih poveznica. Tada je $IxIxK$ tenzor susjedstva \mathcal{X} definiran kao u (3). Primijetio da je tenzor \mathcal{X} rijedak zato što većina stranica pokazuje samo na 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 IxR , i matricu \mathbf{T} veličine KxR , TO YIELD. 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 λ_r definira težinu grupiranja. (Bez smanjenja općenitosti, pretpostavljamo da su stupci naših matrica normalizirani.)

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 $I_1xI_2x\dots xI_n$ i željeni rang aproksimacije R . Tada želimo naći matrice \mathbf{U}^n veličine I_nxR , za svaki $n = 1, \dots, N$ i težinski vektor $\boldsymbol{\lambda}$ duljine R tako da vrijedi:

$$\mathcal{X} \approx \boldsymbol{\lambda}[[\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(N)}]] \quad (5)$$

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

$$\boldsymbol{\lambda}[[\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(N)}]] \equiv \sum_{r=1}^R \lambda_r \mathbf{u}_r^{(1)} \circ \mathbf{u}_r^{(2)} \circ \dots \circ \mathbf{u}_r^{(N)} \quad (6)$$

Bez smanjenja općenitosti, pretpostavljamo da je $\|\mathbf{u}_r^{(n)}\| = 1$ 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:

$$\min \left\| \mathbf{X} - \boldsymbol{\lambda}[[\mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \dots, \mathbf{U}^{(N)}]] \right\|^2 \quad (7)$$

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

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

$$\mathbf{H} \equiv U^{(1)}, \mathbf{A} \equiv U^{(2)}, \mathbf{T} \equiv U^{(3)} \quad (8)$$

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:

$$[\mathbf{X}_{(n)}]_{i,j} = x_{i_1, i_2, \dots, i_N}$$

$$i = i_n, j = 1 + \sum_{k=1, k \neq n}^N (i_k - 1) \prod_{l=1, l \neq k} I_l$$

$$i \leq i \leq I_n, 1 \leq j \leq J \quad (9)$$

3.2.1 Pohlepni PARAFAC - pseudokod

Algorithm 1 Greedy PARAFAC

in: Tensor \mathcal{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 $\{\mathbf{u}_r^{(1)}, \dots, \mathbf{u}_r^{(N)}\}$
 Za svaki $r, r = 1, \dots, N$, unutarnja petlja metodom alternirajućih najmanjih kvadrata minimizira

2.
$$\left\| \left(\mathbf{X} - \sum_{i=1}^{r-1} \lambda_i \mathbf{u}_i^{(1)} \circ \dots \circ \mathbf{u}_i^{(N)} \right) - \left(\mathbf{v}^{(1)} \circ \dots \circ \mathbf{v}^{(N)} \right) \right\| \quad (10)$$

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

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

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:

$$\min_V \left\| \mathbf{X} - [\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(n-1)}, \mathbf{V}, \mathbf{U}^{(n+1)}, \dots, \mathbf{U}^{(N)}] \right\|^2 \quad (11)$$

što se matrično može zapisati kao

$$\min_V \left\| \mathbf{X}_n - \mathbf{V} \mathbf{Z}^{(n)T} \right\|^2 \quad (12)$$

Matrica $\mathbf{Z}^{(n)}$ je veličine $J \times R$ i definirana kao ?? Rješenje najmanjih kvadrata za 12 uključuje pseudo inverz matrice $\mathbf{Z}^{(n)}$:

$$\mathbf{V} = \mathbf{X}_{(n)} (\mathbf{Z}^{(n)T})^\dagger \quad (13)$$

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:

$$(\mathbf{Z}^{(n)T})^\dagger = \mathbf{Z}^{(n)} \mathbf{Y}^{(n)T} \quad (14)$$

Dakle rješenje za 12 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

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

$$\mathbf{X}_{(n)} \mathbf{Z}^{(n)} \tag{15}$$

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

$$\mathcal{X} = \lambda[\mathbf{H}, \mathbf{A}, \mathbf{T}] \quad (16)$$

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 λ_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:

$$q_k = \begin{cases} 1 & \text{ako je term } k \text{ u upitu} \\ 0 & \text{inače} \end{cases} \quad (17)$$

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 \mathbf{s} duljine R koji sadrži vrijednosti svakog grupiranja, na temelju matrice \mathbf{T} iz PARAFAC modela:

$$\mathbf{s} = \mathbf{\Lambda} \mathbf{T}^T \mathbf{q}, \mathbf{\Lambda} = \text{diag}(\boldsymbol{\lambda}) \quad (18)$$

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, $\hat{\mathbf{q}} \in R^l$ te izračunati rezultate grupe kao:

$$\hat{\mathbf{s}} = \mathbf{\Lambda} \mathbf{A}^T \hat{\mathbf{q}}, \mathbf{\Lambda} = \text{diag}(\boldsymbol{\lambda}) \quad (19)$$

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 \mathbf{s} kao u 18. Tada su kombinirani autoriteti dani s:

$$\mathbf{a}^* = \mathbf{A} \mathbf{s} = \sum_{r=1}^R s_r \mathbf{a}_r \quad (20)$$

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

$$\mathbf{h}^* = \mathbf{H} \mathbf{s} = \sum_{r=1}^R s_r \mathbf{h}_r \quad (21)$$

5 Implementacija