

UNIVERZITET U BEOGRADU
MATEMATIČKI FAKULTET



Miloš P. Miković

ALGORITMI ZA REŠAVANJE PROBLEMA
NAJKRAĆE ZAJEDNIČKE NADNISKE

master rad

Beograd, 2021.

Mentor:

dr Aleksandar KARTELJ, docent
Univerzitet u Beogradu, Matematički fakultet

Članovi komisije:

dr Vladimir FILIPOVIĆ, redovni profesor
Univerzitet u Beogradu, Matematički fakultet

dr Stefan MIŠKOVIĆ, docent
Univerzitet u Beogradu, Matematički fakultet

Datum odbrane: _____

Hvala profesoru Aleksandru Kartelju.

Naslov master rada: Algoritmi za rešavanje problema najkraće zajedničke nadni-
ske

Rezime:

Ključne reči: optimizacija, pretraga bima, analiza bioloških sekvenci

Sadržaj

1	Uvod	1
1.1	Problem najkraće zajedničke nadniske	1
1.2	Pregled dosadašnjih istraživanja	2
2	Opis korišćenih optimizacionih metoda	4
2.1	Metaheuristika Pretraga Bima	4
2.2	Heurističke funkcije Većinsko Spajanje i Težinsko Većinsko Spajanje .	6
3	Algoritmi za rešavanje problem najkraće zajedničke nadniske	8
	Bibliografija	9

Glava 1

Uvod

Problem najkraće zajedničke nadniske (*eng.* Shortest Common Supersequence Problem) jedan je od dobro poznatih NP-teških problema optimizacije u oblasti analize reči [1]. Ukratko, PNZN¹ se može opisati kao problem pronalaženja najkraće reči ω sačinjene od simbola zadate konačne Azbuke Σ , tako da su sve sekvence iz unapred zadatog konačnog skupa \mathcal{L} sadržane u sekvenci ω . Kada se kaže da su sve reči iz skupa \mathcal{L} sadržane, misli se na to da se svaka reč iz skupa \mathcal{L} može dobiti uklanjanjem simbola iz reči ω ali u zadatom redosledu [3]. PNZN ima primene u mnogim oblastima informatike uključujući kompresiju podataka [4], optimizaciju upita [13], analizu i poređenje teksta i bioloških sekvenci [15] [2]. Kao rezultat velike primene u mnogim oblastima, postoji veliki broj istraživanja na temu ovog problema u pokušaju da se dođe do što boljeg i prihvatljivijeg rešenja. Trenutno najbolji algoritmi za rešavanje PNZN počivaju na metaheurističkoj metodi pretraga bima (*eng.* beam search) koja će biti predstavljena u ovom radu. U nastavku uvodnog poglavlja biće formalno definisan problem najkraće zajedničke nadniske i biće dat pregled dosadašnjih istraživanja na temu ovog problema.

1.1 Problem najkraće zajedničke nadniske

U ovom poglavlju formalno ćemo definisati PNZN, ali pre toga uvešćemo potrebnu notaciju koja će biti korišćena u nastavku teksta. Konačna azbuka sastoji se od konačnog broja slova i označavaćemo je sa Σ . Svaka konačna reč $\omega = \omega(1)\omega(2)\dots\omega(n)$ sastoji se od konačnog broja slova azbuke gde $\omega(j) \in \Sigma$ predstavlja j -to slovo reči

¹U nastavku teksta PNZN ćemo koristiti kao skraćenicu za problem najkraće zajedničke nadniske

$\omega \in \Sigma^*$. Dužinu reči ω označavaćemo sa $|\omega|$, praznu reč sa ε i važi da $|\varepsilon| = 0$. U skladu sa uvedenom notacijom $|\Sigma|$ predstavlja kardinalnost azbuke. Sa $\omega \supseteq \alpha$ označavaćemo broj pojavljivanja slova α u reči ω ($\omega(1)\omega(2)\dots\omega(n) \supseteq \alpha = \sum_{1 \leq i \leq n, \omega(i)=\alpha} 1$). Reč koja se dobija dodavanjem slova α na početak reči ω označavaćemo sa $\alpha\omega$ (takođe ćemo pisati $\omega = \alpha\omega'$), slično reč koja se dobija skidanjem slova α sa početka reči ω sa $\omega|_\alpha$. Brisanje slova α sa početka svake reči u zadatom skupu, u skladu sa uvedenom notacijom definišemo kao $\{\omega_1, \omega_2, \dots, \omega_n\}|_\alpha = \{\omega_1|_\alpha, \omega_2|_\alpha, \dots, \omega_n|_\alpha\}$.

Neka važi da $\omega_1, \omega_2 \in \Sigma^*$, za reč ω_1 kažemo da je supersekvenca reči ω_2 u oznaci $\omega_1 \succ \omega_2$ ako važi sledeća rekurzivna definicija [1]:

$$\begin{aligned} \omega_1 \succ \varepsilon &\triangleq \text{Tačno} \\ \varepsilon \succ \omega_2 &\triangleq \text{Netačno, Ako } \omega_2 \neq \varepsilon \\ \alpha\omega_1 \succ \alpha\omega_2 &\triangleq \omega_1 \succ \omega_2 \\ \alpha\omega_1 \succ \beta\omega_2 &\triangleq \omega_1 \succ \beta\omega_2, \text{ Ako } \alpha \neq \beta \end{aligned} \tag{1.1}$$

Zapravo, $\omega_1 \succ \omega_2$ označava da se svi simboli iz ω_2 nalaze u ω_1 u datom redosledu, ali ne nužno uzastopno. Na primer, za datu azbuku $\Sigma = \{a, c, t, g\}$, važi $agcatg \succ act$. Sada možemo formalno definisati PNZN. Instanca PNZN može se definisati kao $\mathcal{I} = (\Sigma, \mathcal{L})$, gde Σ predstavlja konačnu azbuku, a \mathcal{L} predstavlja skup od m reči $\{\omega_1, \omega_2, \dots, \omega_m\}$, $\omega_i \in \Sigma^*$. Potrebno je pronaći reč ω najmanje dužine tako da važi da je ω supersekvenca svake reči iz skupa \mathcal{L} ($\omega \succ \omega_i, \forall \omega_i \in \mathcal{L}$ i $|\omega|$ je minimalna). Na primer za instancu PNZN $\mathcal{I} = (\{a, c, t, g\}, \{act, cta, aca\})$, najmanja zajednička supersekvenca instance \mathcal{I} je $acta$.

1.2 Pregled dosadašnjih istraživanja

Problem najkraće zajedničke nadniske prvi je uveo Dejvid Mejer (*eng.* David Maier) 1978. godine u svom radu „The Complexity of Some Problems on Subsequences and Supersequences” [9]. Dokazano je da je PNZN NP-kompletni problem nad svakom azbukom Σ za koju važi da $|\Sigma| \geq 2$ [8]. Korišćenjem dinamičkog programiranja (*eng.* dynamic programming) PNZN nad dve reči dužine n rešen je algoritmom vremenske složenosti $\mathcal{O}(n^2)$ i prostorne složenosti $\mathcal{O}(n^2)$. Algoritam zasnovan na dinamičkom programiranju može biti unapređen, pa tako za k reči dužine maksimalno

n , PNZN može biti rešen u $\mathcal{O}(n^k)$ prostornoj i vremenskoj složenosti [14]. Jasno je da ovakav algoritam nije praktičan za velike vrednosti k . S obzirom na to da ne postoji algoritam polinomijalne složenosti koji rešava PNZN, pribegava se optimizacionim metodama u rešavanju ovog problema. Ono što je karakteristično za optimizacioni pristup rešavanju problema jeste to da se formira algoritam koji rešava postojeći problem tako što daje rešenje koje je prihvatljivo pod određenim uslovima. Takvo rešenje ne mora nužno biti optimalno rešenje problem. Na ovaj način, korišćenjem određene optimizacione tehnike, dobija se algoritam koji se izvršava brzo u realnim uslovima i daje prihvatljivo dobra rešenja.

Vremenom je predloženo mnogo heurističkih i metaheurističkih algoritama za rešavanje PNZN. Neke od poznatijih heurističkih funkcija koje su korišćene u rešavanju PNZN su Alfabet (*eng.* Alphabet) [10], Većinsko Spajanje (*eng.* Majority Merge) i Težinsko Većinsko Spajanje (*eng.* Weighted Majority Merge) [1], Turnirska (*eng.* Tournament) i Pohlepna (*eng.* Greedy) [16], Redukuj-Proširi (*eng.* Reduce-Expand) [10]. Pored navedenih funkcija, korišćeni su i metaheuristički algoritmi, genetski algoritam (*eng.* genetic algorithm) [6] i optimizacija kolonijom (*eng.* colony optimization) [12], koji predstavljaju složenije optimizacione tehnike i imaju tendenciju ka dužem vremenu izvršavanja ne većim instancama problema [7].

Jedna od trenutno najboljih metaheuristika za rešavanje PNZN jeste pretraga bima (*eng.* beam search). Ukratko, pretraga bima predstavlja nepotpunu pretragu stabla, koja na svakom nivou proširuje graf stanja tako što napreduje sa čvorovima koji najviše obećavaju [5]. Upravo zbog toga što se dobro pokazala u rešavanju PNZN i sličnih problema poput problema najduže zajedničke podniske (*eng.* longest common subsequence) pretraga bima će biti korišćena u ovom radu i biće detaljnije opisana u daljem tekstu. S obzirom na to da pretraga bima podrazumeva postojanje heurističke funkcije koja će oceniti kvalitet čvorova na određenom nivou, u ovom radu izabrane su dve takve funkcije, Majority Merge i Weighted Majority Merge koje su se dobro pokazale u prethodnim istraživanjima i one će detaljnije biti opisane u nastavku teksta.

Glava 2

Opis korišćenih optimizacionih metoda

U ovom poglavlju biće dat kratak pregled optimizacionih metoda koje su korišćene za rešavanje problema najkraće zajedničke nadniske. Te metode su:

1. Metaheuristika Pretraga Bima
2. Heuristička funkcija Većinsko Spajanje
3. Heuristička funkcija Težinsko Većinsko Spajanje

U nastavku teksta biće opisane pomenute metode, ali samo kao nezavisne celine, algoritam koji kombinuje rad ovih metoda biće predstavljen u narednom poglavlju.

2.1 Metaheuristika Pretraga Bima

Sada ćemo predstaviti generalnu ideju Pretrage Bima na kojoj ćemo kasnije izgraditi algoritam za PNZN. Pretraga Bima predstavlja metaheuristiku koja je uvedena 1976. godine u oblasti prepoznavanja govora (*eng.* speech recognition). Korišćena je u završnim slojevima mnogih modela za obradu prirodnog jezika (*eng.* natural language processing models) u donošenju odluke da se izabere najbolji izlaz date ciljne promenjive [11]. Pored navedenih primena pretraga bima se intenzivno koristi u sledećim oblastima: kombinatorna optimizacija (*eng.* combinatorial optimization), problemi planiranja (*eng.* scheduling problems), problemi rutiranja vozila (*eng.* vehicle routing problems), podešavanje hiperparametara u oblasti mašinskog učenja



Slika 2.1: Stavi sliku beam searcha

(*eng.* Machine Learning Hyperparameter Tuning), problemi zadovoljenja ograničenja (*eng.* Constraint Satisfaction Problems).

Pretraga Bima predstavlja vrstu pretrage grafa u širinu u cilju da se pronađe najbolji put od korenog čvora do ciljanog čvora u grafu. Kako bi se složenost izračunavanja držala u zadatim granicama, Pretraga Bima evaluira dostignute čvorove na određenom nivou ali bira samo podskup od β čvorova koji najviše obećavaju i sa tim podskupom čvorova napreduje dalje u pretrazi. Izabrani podskup od β čvorova zvaćemo bim (*eng.* beam) i označavaćemo ga sa \mathbb{B} , a β parametar ćemo zvati širina bima (*eng.* beam width) [5]. U kontekstu PNZN graf koji pretražujemo $\mathcal{G} = (\mathcal{V}, \mathbb{E})$ predstavlja usmeren aciklički graf, a pseudo kod algoritma prikazan je pod Algoritam 1.

Algoritam prima tri ulazna parametra: azbuku Σ , ulazni skup reči \mathcal{L} i širinu bima β . Na početku bim \mathcal{B} predstavlja prazan skup. U svakom koraku algoritma bim se

Algoritam 1 PretragaBimaPNZN($\Sigma, \mathcal{L}, \beta$)

```

1:  $\mathcal{B} \leftarrow \{\}$ 
2: while  $True$  do  $\triangleright$  Može postojati uslov izlaska iz petlje pre pronađenog rešenja
3:    $\mathcal{B} \leftarrow \text{ProširiBim}(\Sigma)$ 
4:    $\text{OceniBim}(\mathcal{B}, \mathcal{L})$ 
5:    $\mathcal{B} \leftarrow \text{RedukujBim}(\mathcal{B}, \beta)$ 
6:   if  $\exists \omega \in \mathcal{B}$  tako da  $\text{SuperSekvenca}(\omega, \mathcal{L}) == True$  then
7:     return  $\omega$ 
8:   end if
9: end while

```

proširuje slovima iz azbuke, a zatim se elementi proširenog bima ocenjuje izabranom heurističkom funkcijom. Zatim se vrši odabir β čvorova sa kojima se dalje nastavlja pretraga, i vrši se provera da li neki čvor u bimu predstavlja rešenje problema, ako takav čvor postoji algoritam vraća pronađeno rešenje. Funkcija *ProširiBim* u prvom koraku popunjava bim slovima iz azbuke, a u svakom narednom koraku svako parcijalno rešenje redukovano bima proširuje dodavajući na njega redom slova iz azbuke. Funkcija *OceniBim* predstavlja heurističku funkciju koja dodljuje određene vrednosti svakom elementu bima, a na osnovu ovih vrednosti funkcija *RedukujBim* bira β parcijalnih rešenja koja imaju najveće vrednosti heurističke funkcije. Tokom redukcije bima vrednosti ocena parcijalnih rešenja mogu da se sortiraju tako da redukovani bim sadrži β najbolje ocenjenih rešenja, ali mogu se koristiti i druge strategije odabira rešenja poput ruletske selekcije (*eng.* roulette selection). Algoritam se može zaustaviti kada se pronađe prvo rešenje problema, ali mogu se postaviti i drugi uslovi zaustavljanja poput broja iteracija petlje (predstavlja dubinu unutar stabla pretrage) ili vremenska ograničenja.

2.2 Heurističke funkcije Većinsko Spajanje i Težinsko Većinsko Spajanje

Većinsko Spajanje predstavlja jedan od najpopularnijih algoritama koji rešava PNZN. To je pohlepni algoritam koji inkrementalno konstruiše supersekvencu. Najpre se odredi slovo koje se najčešće nalazi na početku reči iz skupa \mathcal{L} , a zatim se izabrano slovo briše sa početka reči iz \mathcal{L} koje ga sadrže [1]. Pseudo kod algoritma prikazan je pod Algoritam 2.

Mana Većinskog Spajanja je to što ne može da prepozna globalnu strukturu re-

Algoritam 2 VećinskoSpajanje($\mathcal{L} = \{\omega_1, \omega_2, \dots, \omega_m\}, \Sigma$)

```

1:  $\mathcal{S} \leftarrow \varepsilon$  ▷ supersekvencu
2: while  $\sum_{\omega_i \in \mathcal{L}} |\omega_i| \neq 0$  do
3:   for  $\alpha \in \Sigma$  do
4:      $v(\alpha|\mathcal{S}) \leftarrow \sum_{\omega_i \in \mathcal{L}, \omega_i = \alpha\omega'_i} 1$ 
5:   end for
6:    $\beta \leftarrow \operatorname{argmax}\{v(\alpha|\mathcal{S}) \mid \alpha \in \Sigma\}$ 
7:    $\mathcal{L} \leftarrow \mathcal{L}|_\beta$ 
8:    $\mathcal{S} \leftarrow \mathcal{S}\beta$ 
9: end while
10: return  $\mathcal{S}$ 

```

či iz skupa \mathcal{L} . U principu Većinsko spajanje izostavlja činjenicu da reči mogu biti različitih dužina. To dalje znači da će slova sa početka kraćih reči imati veću šansu da budu uklonjena iako algoritam i dalje treba da obradi preostale dugačke reči. Iz tog razloga bi skidanje slova sa početka kraćih reči trebalo da bude manje prioritarno. Drugim rečima bolje je da se prioritizira skidanje slova sa početka dužih reči. To se može postići dodavanjem težine svakom simbolu azbuke u odnosu na dužinu preostalih reči iz \mathcal{L} nakon uklanjanja tog simbola sa početka reči koje počinju tim simbolom. Dakle korak 4 u algoritmu VećinskoSpajanje možemo zameniti sa:

$$v(\alpha|\mathcal{S}) \leftarrow \sum_{\omega_i \in \mathcal{L}, \omega_i = \alpha\omega'_i} |\omega'_i| \quad (2.1)$$

Ovako modifikovan algoritam naziva se Težinsko Većinsko Spajanje i postoje indikacije da na određenim instancama problema može da nadmaši algoritam Većinskog Spajanja, pogotovo kada nema struktuiranosti unutar skupa \mathcal{L} ili kada je ta struktuiranost haotična [1]. Predstavljena dva algoritma biće korišćena u nastavku rada kao heurističke funkcije koje će ocenjivati elemnte bima.

Dodaj tekst za LAWMM...

Glava 3

Algoritmi za rešavanje problem najkraće zajedničke nadniske

dsadsadasd

Bibliografija

- [1] Antonio J. Fernandez Christian Blum, Carlos Cotta and Francisco Gallardo. A Probabilistic Beam Search Approach to the Shortest Common Supersequence Problem. In *Lecture Notes in Computer Science*, 2007.
- [2] Joseph Kruskal David Sankoff. *Time Warps, String Edits, and Macromolecules: The Theory and Practice of Sequence Comparison*. Center for the Study of Language and Inf, 1983.
- [3] Garey and Johnson. Shortest common supersequence. on-line at: <https://www.csc.kth.se/~viggo/wwwcompendium/node165.html>.
- [4] Storer JA. *Data compression: methods and theory*. Computer Science Press, 1988.
- [5] Marc Huber Gunther Raidl Jonas Mayerhofer, Markus Kirchweger. A Beam Search for the Shortest Common Supersequence Problem Guided by an Approximate Expected Length Calculation. 2022.
- [6] Frerk Schneider Jurgen Branke, Martin Middendorf. Improved heuristics and a genetic algorithm for finding short supersequences, 1998.
- [7] Hon Wai Leong Kang Ning. Towards a better solution to the shortest common supersequence problem: the deposition and reduction algorithm, 2006.
- [8] Esko Ukkonen Kari-Jouko Raiha. The shortest common supersequence problem over binary alphabet is NP-complete. In *Theoretical Computer Science*, 1981. Volume 16, Issue 2, Pages 187-198.
- [9] David Maier. The Complexity of Some Problems on Subsequences and Supersequences, 1978.

- [10] Gianluca Della Vedova Giancarlo Mauri Paolo Barone, Paola Bonizzoni. An Approximation Algorithm for the Shortest Common Supersequence Problem: An Experimental Analysis, 2001.
- [11] Matt Payne. What is Beam Search? Explaining The Beam Search Algorithm. on-line at: <https://www.width.ai/post/what-is-beam-search>.
- [12] Martin Middendorf Rene Michel. An island model based ant system with lookahead for the shortest supersequence problem. 2006.
- [13] Timos K. Sellis. Multiple-query optimization, 1988. ACM Transactions on Database Systems (TODS), 13(1):23-52.
- [14] Ming Li Tao Jiang. On the Approximation of Shortest Common Supersequences and Longest Common Subsequences, 1995.
- [15] Ronald L. Rivest Clifford Stein Thomas H. Cormen, Charles E. Leiserson. *Introduction to Algorithms, Second edition*. MIT Press and McGraw-Hill, 2001.
- [16] Vadim G. Timkovsky. Some Approximations for Shortest Common Nonsubsequences and Supersequences. In *String Processing and Information Retrieval*, 2006.

Biografija autora

Vuk Stefanović Karadžić (*Tršić, 26. oktobar/6. novembar 1787. — Beč, 7. februar 1864.*) bio je srpski filolog, reformator srpskog jezika, sakupljač narodnih umotvorina i pisac prvog rečnika srpskog jezika. Vuk je najznačajnija ličnost srpske književnosti prve polovine XIX veka. Stekao je i nekoliko počasnih mastera. Učestvovao je u Prvom srpskom ustanku kao pisar i činovnik u Negotinskoj krajini, a nakon sloma ustanka preselio se u Beč, 1813. godine. Tu je upoznao Jerneja Kopitara, cenzora slovenskih knjiga, na čiji je podsticaj krenuo u prikupljanje srpskih narodnih pesama, reformu ćirilice i borbu za uvođenje narodnog jezika u srpsku književnost. Vukovim reformama u srpski jezik je uveden fonetski pravopis, a srpski jezik je potisnuo slavenosrpski jezik koji je u to vreme bio jezik obrazovanih ljudi. Tako se kao najvažnije godine Vukove reforme ističu 1818., 1836., 1839., 1847. i 1852.