

Univerzitet u Beogradu
Matematički fakultet

Seminarski rad

Analiza, klasifikacija i klasterovanje proteinskih blokova

Mentor:

Prof. dr Nenad Mitić
Katedra za računarstvo i informatiku

Studenti:

Anja Milutinović 235/2021
Đurđa Milošević 84/2021
Smer: Informatika

Datum: 2024/25

Sadržaj

1	Uvod	2
2	O proteinskim blokovima	3
3	Analiza	4
3.1	Opis podataka	4
3.1.1	Izdvajanje retkih prelaza između proteinskih blokova u podacima	6
3.2	Izdvajanje parova aminokiselina prisutnih u neočekivanim prelazima	6
3.3	Izdvajanje parova sekundarnih struktura koji se javljaju u neočekivanim prelazima	8
3.4	Analiza vrednosti pLDDT parametra	8
3.5	Analiza vrednosti RSA parametra	8
3.6	Ispitivanje zastupljenosti aminokiselina u podacima	8
4	Klasifikacija	9
5	Klasterovanje	10
6	Zaključak	11

1 Uvod

Proteini su osnovni gradivni blokovi svih ćelija u organizmu i kao takvi igraju ključnu ulogu u održavanju života, reprodukciji, odbrani i replikaciji. Sve funkcije proteina zavise od njihove strukture, zbog čega je analiza strukture proteina od izuzetnog značaja u bioinformatičkim i biohemijskim istraživanjima.

U nastojanju da se postigne precizniji, detaljniji i informativniji prikaz trodimenzionalne strukture proteina, razvijeni su proteinski blokovi. Često nazvani i strukturnim alfabedom, pokazano je da omogućavaju detekciju strukturne sličnosti između proteina sa izuzetnom efikasnošću.

Ovaj seminarski rad se fokusira na proteinske blokove, sa posebnim akcentom na retke i neočekivane prelaze između njih. Niz proteinskih blokova dobijen je analizom humanog proteoma, koji je generisan pomoću AlphaFold2 programa. Seminarski rad se sastoji od tri ključna segmenta: analize, klasifikacije i klasterovanja proteinskih blokova.

Analiza obuhvata izdvajanje aminokiselina i sekundarnih struktura u retkim prelazima, određivanje prirode vrednosti pLDDT (*the predicted local distance difference test*) i RSA (*relative solvent accessibility*) parametara, kao i poređenje zastupljenosti pojedinačnih aminokiselina u podacima u odnosu na očekivane procenete.

Klasifikacija je usredsređena na predviđanje aminokiselina i sekundarnih struktura u prelazima, kao i parova proteinskih blokova koji čine posmatrane prelaze. Ovaj pristup omogućava donošenje interesantnih zaključaka o relevantnosti, ubedljivosti i stabilnosti proteinskih blokova kao strukturnih indikatora.

Klasterovanje je izvedeno nad skupom podataka koji sadrži informacije o strukturi proteina, uključujući prelaze između proteinskih blokova, kao i aminokiseline i sekundarne strukture prisutne u tim prelazima.

Seminarski rad je realizovan u okviru kursa „Istraživanje podataka 2” na Matematičkom fakultetu Univerziteta u Beogradu.

2 O proteinskim blokovima

3 Analiza

U ovom poglavlju analizirani su prelazi između proteinskih blokova. Fokus analize bio je na identifikaciji neočekivanih prelaza između proteinskih blokova i proučavanju aminokiselina i sekundarnim struktura prisutnih u njima. Dodatno, ispitane su vrednosti relevantnih parametara, kao i zastupljenost pojedinačnih aminokiselina u podacima.

3.1 Opis podataka

Podaci korišćeni u analizi obuhvataju informacije o proteinskim blokovima (PBs), aminokiselinama (AA), sekundarnim strukturama (S2), predviđenoj učestalosti prelaza, kao i dodatnim parametrima poput pLDDT i RSA. Izvor podataka su rezultati generisani pomoću AlphaFold2 programa. Konkretno, analiza je sprovedena nad dve datoteke: prvobitnog, obimnijeg skupa podataka, koji sadrži sledeće informacije:

```
Protein_number  res_number  PB1  PB2  AA1  AA2  S2_1  S2_2
0 7 j j G G C S
expected_frequency  pLDDT  RSA1  RSA2
0.880055  60.180000  100.000000  100.000000
```

i manjeg podskupa dobijenog filtriranjem rezultata u skladu sa humanim proteomom, čija je struktura sledeća:

```
Protein_number PB1 PB2 expected_frequency AA1 AA2 S2_1 S2_2
1 p f 0.94129547 H E C S
```

Podaci su organizovani u tabelarnom formatu, pri čemu svaki red sadrži informacije o konkretnom prelazu između proteinskih blokova, aminokiselinama uključenim u taj prelaz, pripadajućim sekundarnim strukturama, predviđenoj učestalosti prelaza i vrednostima parametara pLDDT i RSA.

U nastavku je dat detaljan opis atributa:

- **Protein_number** – diskretan, kategorijski i redni atribut koji označava pojedinačne proteine u skupu podataka.
- **res_number** – diskretan, kategorijski i redni atribut koji označava poziciju aminokiselinskog ostatka unutar sekvence.
- **PB1, PB2** – diskretni, kategorijski i nominalni atributi koji predstavljaju oznake proteinskih blokova između kojih dolazi do prelaza.

- **AA1, AA2** – diskretni, kategorijski i nominalni atributi koji opisuju aminokiseline prisutne u prelazu. Postoji ukupno 20 različitih vrednosti za ove attribute, u skladu sa standardnim skupom aminokiselina koje grade proteine. U tabeli 1 prikazane su odgovarajuće jednoslovne oznake za svaku aminokiselinu.

Aminokiselina	Jednoslovna oznaka
Alanin	A
Arginin	R
Asparagin	N
Asparaginska kiselina	D
Cistein	C
Glutaminska kiselina	E
Glutamin	Q
Glicin	G
Histidin	H
Izoleucin	I
Leucin	L
Lizin	K
Metionin	M
Fenilalanin	F
Prolin	P
Serin	S
Treonin	T
Triptofan	W
Tirozin	Y
Valin	V

Tabela 1: Aminokiseline i njihove jednoslovne oznake korišćene kao vrednosti atributa AA1 i AA2.

- **S2_1, S2_2** – diskretni, kategorijski i nominalni atributi koji predstavljaju sekundarne strukture prisutne u prelazu. Vrednosti ovih atributa određene su u skladu sa DSSP (*Dictionary of Secondary Structure of Proteins*) standardom, koji klasifikuje sekundarne strukture na osnovu vodoničnih veza između amino i karboksilnih grupa i geometrijskih karakteristika polipeptidnog lanca. U tabeli 2 prikazane su oznake i odgovarajući tipovi sekundarne strukture [1].

Oznaka	Tip sekundarne strukture
H	Alfa-heliks (α - <i>helix</i>)
B	Beta-most (<i>beta bridge</i>)
E	Prošireni beta list (<i>extended beta sheet</i>)
G	3(10)-heliks
I	Pi-heliks (π - <i>heliks</i>)
T	Heliks okret (<i>helix-turn</i>)
S	Zavojnica (<i>bend</i>)
L	Nespecifična ili neorganizovana struktura (<i>loop</i>)

Tabela 2: Sekundarne strukture proteina prema DSSP standardu, korišćene kao vrednosti atributa S2_1 i S2_2.

- **expected_frequency** - neprekidan, kvantitativan i razmerni atribut čija vrednost pripada intervalu $[0,1]$ i označava očekivanu učestalost prelaza između proteinskih blokova.
- **pLDDT, RSA1, RSA2** – neprekidni, kvantitativni i razmerni atributi čije su vrednosti u intervalu $[0,100]$. Parametar pLDDT procenjuje pouzdanost predikcije lokalne strukture proteina, dok RSA1 i RSA2 predstavljaju relativnu dostupnost aminokiselinskih ostataka rastvaraču.

3.1.1 Izdvajanje retkih prelaza između proteinskih blokova u podacima

Glavni zadatak analize bio je proučavanje prelaza između proteinskih blokova pri čemu je akcenat stavljen na neočekivane, retke prelaze. Takvi prelazi su, zbog svoje prirode, od posebnog interesa za istraživanje jer mogu doprineti novim saznanjima o nizu proteinskih blokova koji opisuje trodimenzionalnu strukturu proteina.

U ovom radu, retkim prelazom smatra se onaj koji se javlja u manje od 1% slučajeva.

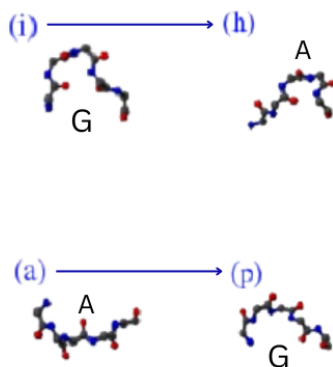
3.2 Izdvajanje parova aminokiselina prisutnih u neočekivanim prelazima

U cilju pronalaženja potencijalnih korelacija između atributa, odnosno između parova aminokiselina vezanih za retki prelaz i proteinskih blokova koji formiraju prelaz, izvršeno je izdvajanje tih parova. Pored toga, izračunate su frekvencije svih parova, a zatim su izdvojeni oni najfrekventiniji.

S obzirom na to da se u podacima korišćenim za ovu analizu prelazi ne nado-
vezuju, posmatran je pojedinačno svaki prelaz i unutar njega određen svaki
par aminokiselina. Važno je napomenuti da redosled aminokiselina u paru ni-
je irelevantan, zbog čega se parovi aminokiselina (A1, A2) ne mogu smatrati
identičnim parovima (A2, A1). Kako bi ova tvrdnja bila intuitivnija i jasnija,
potrebno je prvo objasniti zašto je prelaz određen sa dve aminokiseline iako
su proteinski blokovi sačinjeni od 5 uzastopnih aminokiselina.

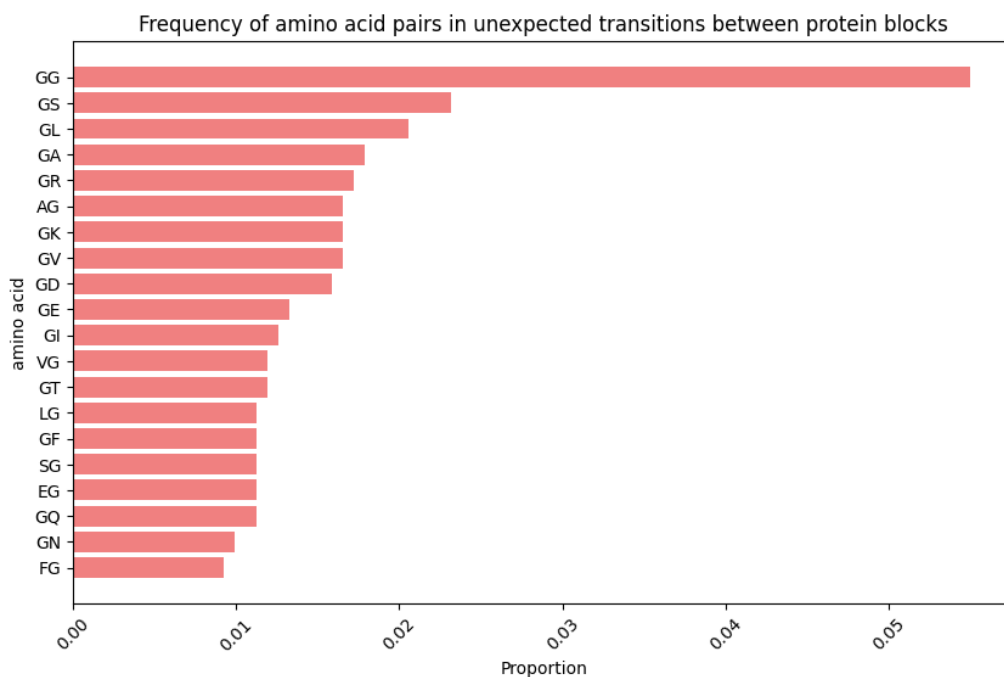
Naime, prilikom dodeljivanja proteinskog bloka fragmentu od 5 aminokise-
lina, centralna aminokiselina igra ključnu ulogu jer je ona ta koja efektivno
definiše fragment i omogućava izračunavanje 8 diedarskih uglova $\psi_{i-2}, \phi_{i-1},$
 $\psi_{i-1}, \phi_i, \psi_i, \phi_{i+1}, \psi_{i+1}, \phi_{i+2}$. Ovi uglovi se zatim upoređuju sa već definisanim
diedarskim uglovima prototipova proteinskih blokova nakon čega se dodeljuje
odgovarajući proteinski blok. Dakle, za svaki prelaz centralna aminokiselina
je najvažnija zbog njenog uticaja na ceo fragment i zato je baš ona ta koja
je izabrana da bude deo podataka, odnosno da bude predstavnik u prelazu.

Imajući u vidu prethodno navedeno, obrnut redosled aminokiselina u paru
može da predstavlja potpuno različit prelaz između proteinskih blokova i zato
je važno analizirati ih kao zasebne parove. Slika 1 prikazuje primer iz skupa
podataka koji ilustruje iznesenu tvrdnju. Izračunate su frekvencije svakog



Slika 1: Redosled aminokiselina u paru je važan.

para aminokiselina u prelazima. Najčešće se javlja par (G, G). Pregled ostalih
19 najučestalijih parova dat je na slici 2.



Slika 2: Frekvencija parova aminokiselina u retkim prelazima.

Dobijeni najfrekventniji par je od posebnog interesa s obzirom na specifične karakteristike glicina. Glicin ima najvažniju ulogu u formiranju sekundarne strukture alfa heliksa zahvaljujući svojoj fleksibilnosti, koja je rezultat izuzetno malog bočnog lanca. Takođe, glicin doprinosi stabilnosti tercijalne strukture proteina [2]. Na osnovu ovoga, može se zaključiti da se u većini retkih prelaza postiže strukturna stabilnost zahvaljujući prisustvu glicina.

3.3 Izdvajanje parova sekundarnih struktura koji se javljaju u neočekivanim prelazima

3.4 Analiza vrednosti pLDDT parametra

3.5 Analiza vrednosti RSA parametra

3.6 Ispitivanje zastupljenosti aminokiselina u podacima

4 Klasifikacija

5 Klasterovanje

6 Zaključak

Literatura

- [1] Phil Carter, Claus A. F. Andersen, Burkhard Rost. *DSSPcont: Continuous Secondary Structure Assignments for Proteins*. Bioinformatics, 19(2):230-231, 2003.
- [2] Hao Dong, Mukesh Sharma, Huan-Xiang Zhou, Timothy A Cross. *Glycines: Role in α -Helical Membrane Protein Structures and a Potential Indicator for Native Conformation*. Biophysical Journal, 109(1):1-8, 2015.