

# Playing Regex Golf with Genetic Programming

Projekat u okviru kursa Računarska inteligencija  
Matematički fakultet  
Univerzitet u Beogradu

Anđela Ilić  
[mi17105@alas.matf.bg.ac.rs](mailto:mi17105@alas.matf.bg.ac.rs)  
Mina Milošević  
[mi17081@alas.matf.bg.ac.rs](mailto:mi17081@alas.matf.bg.ac.rs)

Februar 2021

# Sadržaj

<b>1</b>	<b>Opis problema</b>	<b>3</b>
<b>2</b>	<b>Implementacija</b>	<b>3</b>
2.1	Obrada ulaznih podataka . . . . .	3
2.2	Genetsko programiranje . . . . .	4
2.2.1	Implementacija jedinki . . . . .	4
2.2.2	Parametri genetskog programiranja . . . . .	4
<b>3</b>	<b>Rezultati</b>	<b>4</b>
<b>4</b>	<b>Zaključak</b>	<b>4</b>
<b>5</b>	<b>Izvori</b>	<b>4</b>

## 1 Opis problema

Data su dva skupa reči -  $M$  i  $U$ . Cilj *Regex Golf* igre je pronaći najkraći regularni izraz kojim se mogu zapisati sve reči iz skupa  $M$ , ali kojim se ne može zapisati nijedna reč skupa  $U$ . Za date skupove  $M$  i  $U$  ne možemo sa sigurnošću da tvrdimo da postoji rešenje koje zadovoljava prethodne uslove. Takođe, ako dobijemo regularni izraz koji zadovoljava navedene uslove, ne možemo za svaki primer znati da li postoji i bolje rešenje tj. kraći regularni izraz.

## 2 Implementacija

Svaka jedinka u Genetskom programiranju će biti predstavljena kao drvo. U listovima nalaze elementi koje ćemo jednim imenom zvati *Terminali* (terminal set), a u unutrašnjim čvorovima su elementi koje nazivamo *Funkcije* (function set).

Skup funkcija sadrži operatore koji se mogu javiti u regularnim izrazima. Primeri takvih operatora su: `.*`, `++`, `.*+`, `..`, `.*`, `(.)`, `[]`, `[ ]`, `..`, `..|`. Tačka `.` je mesto na kome se nalaze deca u drvetu.

Skup terminala čine elementi koji zavise i koji ne zavise od ulaznih skupova  $M$  i  $U$ . Elementi koji su nezavisni - opsezi malih i velikih slova, brojeva u regularnim izrazima, karakteri  $\wedge$  i  $\$$ , wildcard karakter  $\cdot$ . Elementi skupa terminala koji su zavisni - skup karaktera iz  $M$ , parcijalni opsezi karaktera iz  $M$  i  $n$ -grami.

## 2.1 Obrada ulaznih podataka

Bez obzira na veličinu ulaznih skupova  $M$  i  $U$  i njihov sadržaj, implementacija prethodno navedenih zavisnih terminala je ista.

*Skup karaktera iz  $M$  sadrži sve karaktere koji se mogu naći u rečima iz  $M$ , bez ponavljanja i sortirani po engleskom alfabetu.*

*Parcijalni opsezi* se prave na osnovu skupa karaktera iz  $M$ . Potrebno je naći maksimalne podskupove tog skupa karaktera tako da se svi karakteri iz intervala  $[c_f, c_l]$  nalaze u skupu karaktera iz  $M$ .  $c_f$  je prvi karakter, a  $c_l$  je poslednji karakter iz podskupa. Kao rezultat se vraćaju parcijalni opsezi u formatu  $c_f - c_l$ .

*n*-grami. Pravimo skup svih *n*-grama dužine  $2 \leq n \leq 4$  koji se mogu naći u *M* ili u *U* (ili oba). Svakom dobijenom *n*-gramu se dodelju vrednost koja predstavlja njegov *score*. Za svaku reč iz *M* koja sadrži dati *n*-gram, njegov *score* se uvećava za 1, a za svaku reč iz *U* koja ga sadrži, *score* se umanjuje za 1. Nakon formiranja svih *n*-grama i određivanja njihovih vrednosti, sortiraju se opadajuće po vrednosti. Potrebno je uzeti najmanji podskup *n*-grama tako da njihova ukupna vrednost (*score*) bude jednaka bar dužini skupa *M* ( $|M|$ ).

## **2.2 Genetsko programiranje**

### **2.2.1 Implementacija jedinki**

### **2.2.2 Parametri genetskog programiranja**

## **3 Rezultati**

## **4 Zaključak**

## **5 Izvori**