

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

ZAVRŠNI RAD br. 4334

Primjena genetskog programiranja na problem klasifikacije podataka

Ivan Vlašić

Zagreb, lipanj 2016.

Umjesto ove stranice umetnite izvornik Vašeg rada.
Da bi ste uklonili ovu stranicu obrišite naredbu \izvornik.

SADRŽAJ

1. Uvod	1
2. Strojno učenje i problem klasifikacije	2
2.1. Strojno učenje	2
2.2. Klasifikacija podataka	3
2.3. Najznačajniji algoritmi za rješavanje problema klasifikacije	4
3. Genetsko programiranje	6
3.1. Stvaranje početne populacije	7
3.1.1. Metoda full	7
3.1.2. Metoda grow	7
3.1.3. Metoda ramped half-and-half	7
3.2. Procjena dobrote jedinke	7
3.3. Genetski operatori	8
3.3.1. Selekcija	8
3.3.2. Križanje	9
3.3.3. Mutacija	9
3.4. Problemi genetskog programiranja	10
3.5. Primjena genetskog programiranja na problem klasifikacije	11
4. Programsko ostvarenje	13
4.1. Radno okruženje ECF	13
4.2. Korišteni klasifikatori	14
4.3. Funkcije čvorova	15
4.4. Funkcije dobrote	16
4.5. Parametri i konfiguracijske datoteke	17
5. Primjena i analiza rezultata	20
5.1. <i>Iris flower data set</i>	20

5.1.1. Klasifikacija stakla	22
6. Zaključak	26
Literatura	27

1. Uvod

Klasifikacija je problem razvrstavanja podataka u kategorije na osnovu skupine podataka za učenje čija je podjela već poznata. Kao primjer koji se u praksi često koristi, možemo navesti prepoznavanje spada li novi email pod neželjeni sadržaj ili ne. Razvojem računarske znanosti, a pogotovo strojnog učenja, došlo je do pojave različitih metoda rješavanja tog problema. Neke su bile više uspješne od drugih. Ovaj rad će pokušati pojasniti rješavanje problema klasifikacije podataka upotrebom genetskog programiranja koje se pokazalo kao dosta precizna metoda iako može dovesti do problema prevelike složenosti kod velike skupine podataka.

U početku ćemo se upoznati s problemom klasifikacije, njegovim najčešćim primjenama i prihvaćenim metodama rješavanja te osnovnim karakteristikama genetskog programiranja koje ga čine pogodnim za rješavanje ovog problema. Posebnu pažnju ćemo posvetiti algoritmima klasifikacije u obliku regresijskog modela i logičkih pravila. Pojasnit ćemo programski sustav implementiran uz ovaj rad, modele i operatore koje smo koristili kao i detaljne upute za korištenje. Analizirat ćemo rezultate dobivene primjenom navedenog sustava na testne primjere te donijeti zaključke o njegovoj optimiziranosti i kvaliteti dobivenih rješenja.

2. Strojno učenje i problem klasifikacije

U ovom poglavlju ćemo dati kratki pregled strojnog učenja i upoznati se s njegovim vrstama. Objasnit ćemo problem klasifikacije podataka i proučiti njegove najčešće primjene. Na kraju ćemo opisati neke od najznačajnijih algoritama koji se u današnje vrijeme koriste za rješavanje tog problema kao i njihovu uspješnost.

2.1. Strojno učenje

Strojno učenje je grana računarske znanosti koja se razvila iz prepoznavanja uzoraka i teorije učenja u umjetnoj inteligenciji. Podrazumijeva programiranje računala tako da optimiziraju neki kriterij uspješnosti temeljem podatkovnih primjera ili prethodnog iskustva. Strojno učenje je veliki razvoj doživjelo osamdesetih godina prošlog stoljeća kad se počelo udaljavati od umjetne inteligencije i razvijati kao zasebna grana računarske znanosti.

Strojno učenje možemo podijeliti na nadzirano, nenadzirano i podržano učenje. Nadzirano učenje se koristi ulaznim skupom podataka kako bi učilo i izgradilo model predviđanja. Podaci za učenje se sastoje od primjera koji imaju vektor ulaznih podataka i predviđen izlaz. Algoritmima nadziranog učenja se ti podaci analiziraju i stvara se funkcija koja se može koristiti za mapiranje novih primjera. Ovim postupcima mogu se rješavati dvije vrste problema: klasifikacija i regresija. Kod klasifikacije primjeru pridružujemo klasu kojoj taj primjer pripada, dok kod regresije primjeru pridružujemo neku kontinuiranu vrijednost.

Za razliku od nadziranog učenja, nenadzirano učenje ne koristi podatke za učenje. Kod ove vrste učenja dani su nam podaci bez ciljne vrijednosti i naš je zadatak pronaći pravilnosti među njima. Tipično se koristi kod eksplorativne dubinske analize podataka, u biologiji za grupiranje organizama prema njihovim značajkama, grupiranje sličnih dokumenata, grupiranje DNA-mikropolja i slično. Koriste se grupiranje,

procjena gustoće i smanjenje dimenzionalnosti.

Podržano učenje je učenje strategije na temelju serije izlaza. Tipično se koristi kod igara, robotike i upravljanja te višeagentnih sustavima. Pristup temelji na agentima koji ovisno o kvaliteti postignutih rezultata postižu nagrade. Cilj svakog agenta je maksimizirati dobivenu nagradu i time temelji uspjeh svojih strategija.

2.2. Klasifikacija podataka

Kao što je već objašnjeno u uvodu poglavlja, problemi klasifikacije podataka se rješavaju postupcima nadziranog učenja. Svakom primjeru pridružujemo određenu klasu, odnosno razred kojoj pripada. Primjer možemo definirati kao vektor značajki gdje n predstavlja dimenziju vektora, $X = (x_1, x_2, \dots, x_n)^T$. Vektor predstavlja ulazni prostor (engl. *input space*) ili prostor primjera (engl. *instance space*). Kod nadziranog učenja unaprijed nam je poznata klasa kojoj pripada primjer iz skupa učenja. Cilj klasifikacije je određivanje nepoznatih klasa ili grupa u podacima, to jest svrstati nove, već neviđene podatke u neki od već unaprijed poznatih klasa. Najjednostavniji je slučaj u kojem se podaci moraju svrstati u samo jedan razred. Kod više razreda, radi se o klasifikaciji s višestrukim oznakama. Ako se ograničimo samo na dvije klase, klasifikator nazivamo binarni klasifikator.

Skup primjera za učenje sastoji se od parova primjera i pripadnih oznaka. Možemo ga prikazati tablično kao:

x_1	x_2	...	y
$x_1^{(1)}$	$x_2^{(1)}$...	$y^{(1)}$
$x_1^{(2)}$	$x_2^{(2)}$...	$y^{(1)}$
\vdots	\vdots		\vdots
$x_N^{(1)}$	$x_N^{(1)}$...	$y^{(N)}$

gdje je N ukupan broj primjera, a i indeks primjera. Skup ulaznih oznaka označavamo s $(x^{(i)})_{i=1}^N$, a izlaznih s y . Zadaća klasifikacijskog algoritma jest naučiti hipotezu koja određuje pripada li neki primjer određenoj klasi ili ne. Da bi u tome uspio, potrebno je izgraditi klasifikator. Ako klasifikator dobro generalizira ispravno će klasificirati još neviđene podatke. Važno je napomenuti da uspjeh klasifikatora u velikoj mjeri ovisi o svojstvima podataka koje je potrebno klasificirati. Ne postoji klasifikator koji radi najbolje za sve probleme (pojava koja je poznatija pod imenom *no free*

lunch theorem). Određivanje primjerenog klasifikatora za neki problem je dosta složen i težak problem koji nije uvijek moguće riješiti. Učenje hipoteze je loše definiran problem (engl. *ill-posed problem*). Primjeri za učenje nisu sami po sebi dovoljni da bi se hipoteza jednoznačno definirala. Svojstvo hipoteze da odredi klase još neviđenih primjera naziva se generalizacija.

Jedan od glavnih problema do kojih može doći je problem šuma. Šum je neželjena anomalija u podacima koja je uzrokovana različitim uzrocima poput nepreciznosti pri mjerenju značajki, pogreški u označavanju, nejasnih granica klasa i postojanju skrivenih značajki. U načelu, šum nije moguće razdvojiti od pravih podataka. Danas klasifikacija ima mnoge primjene. Koristi se u računalnom vidu, prepoznavanju govora, klasifikaciji dokumenata i slično. Svi ovi i mnogi drugi problemi se rješavaju različitim algoritmima od kojih svi imaju svoje pozitivne i negativne strane, kao i različite uspješnosti. U sljedećem poglavlju ćemo razmotriti neke od najčešćih algoritama koji se koriste u praksi kao i njihove značajke.

2.3. Najznačajniji algoritmi za rješavanje problema klasifikacije

Algoritama za rješavanje problema klasifikacije postoji mnogo. Razlikuju se po svojoj jednostavnosti primjene, točnosti i razumljivosti. Uspješnost algoritama u velikoj mjeri ovisi o problemu kojeg rješavamo. Ponekad nam je važnija brzina klasifikatora od njegove točnosti, a u drugim slučajevima je važnije da klasifikator ima veliku točnost nego brzinu ili jednostavnost. Među algoritme linearne klasifikacije spadaju Fisherova linearna diskriminanta i Naivni Bayes klasifikatori koji su se u početku pokazali dosta učinkovitim za rješavanje složenih problema, no detaljna analiza 2006. godine je pokazala da postoje drugi algoritmi koji su dosta bolji, poput algoritma *random forest*. U tu skupinu još spadaju logistička regresija i *perceptron*. Jedan od češće korištenih algoritama je i metoda potpornih vektora (engl. *Support Vector Machine*) koji se bazira na učenju razdvajanjem funkcija i estimacije funkcije u regresiji. U velikoj upotrebi su i neuronske mreže koje su u širom smislu replika ljudskog mozga kojom se nastoji simulirati postupak učenja. To je skup međusobno povezanih jednostavnih procesnih elemenata koji se nazivaju čvorovi ili neuroni, i čija se funkcionalnost zasniva na biološkom neuronu. Neuronske mreže su vrlo dobre u procjeni nelinearnih odnosa uzoraka, mogu raditi s nejasnim ili manjkavim podacima, mogu raditi s velikim brojem varijabli i prilagođavati se okolini što ih čini pogodnim za probleme klasifikacije. Još neki

od poznatijih metoda i algoritama su kvadratna metoda koja je dizajnirana za situacije kada karakteristike svake grupe imaju normalnu distribuciju te metoda kernel procjene.

Za probleme klasifikacije također se može koristiti i genetsko programiranje koje se zasniva na razvoju stabala odluke. Prednost genetskog programiranja leži u činjenici da može otkriti manje vidljive veze između podataka. Važna primjena je postignuta u otkrivanju gena. Više o tome će biti objašnjeno u sljedećih nekoliko poglavlja.

3. Genetsko programiranje

U ovom poglavlju ćemo se detaljnije upoznati s genetskim programiranjem i vrstama koje se danas koriste. Objasnit ćemo probleme za čije je rješavanje pogodno genetsko programiranje kao i načini do kojih algoritmi dolaze do rješenja.

Genetsko programiranje je optimizacijska tehnika koja spada pod skupinu evulucijskih algoritama. Radi na principu simulacije evolucije, bolja rješenja preživljavaju i svoja svojstva prenose u sljedeće generacije dok loša izumiru. Kao i u prirodi, koriste se mehanizmi selekcije, križanja i mutacije čime možemo doći do raznolikog broja jedinki i uspješno se približiti dovoljno dobrom rješenju promatranog problema. Razvoj genetskog programiranja možemo pratiti od pedesetih godina prošlog stoljeća, no tek šezdesetih i sedamdesetih godina su postali šire prepoznati kao pogodna tehnika za rješavanje optimizacijskih problema. Većem napretku je pridonio John R. Koza koji je GP primijenio na velikom broju optimizacijskih problema i problema pretraživanja različite složenosti. Devedesetih se genetsko programiranje uglavnom koristilo za relativno jednostavne probleme budući da su izračuni dosta računalno intenzivni čemu računala tog doba nisu bila dorasla. U zadnjim godinama, zbog velikog napretka u razvoju procesora, genetsko programiranje se koristi za sve širi spektar problema poput sortiranja, pretraživanja, elektroničkog dizajna i mnogih drugih.

Genetsko programiranje se često uspoređuje sa genetskim algoritmima. Međutim, postoje određene razlike. Kod genetskog programiranja rješenje je prikazano računalnim programom. Rješenje je onaj program koji daje rješenje zadanog problema u u najkraćem roku. Kod genetskih algoritama postoje različite reprezentacije rješenja ovisno o problemu koji rješavamo. Rješenje se može predstavljati brojkom, nizom bitova ili brojeva i drugima. Još jedna od većih razlika je duljina rješenja. Kod genetskog algoritma duljina rješenja se kroz evaluaciju ne mijenja, dok kod GP-a stabla često mijenjaju dubinu, veličinu i širinu. To često dovodi do problema prekomjernog rasta zbog čega je potrebno implementirati određena ograničenja duljine rješenja, poput dubine stabla ili broja čvorova.

U sljedećim potpoglavljima ćemo se detaljnije upoznati o načinu na koji genetsko

programiranje radi.

3.1. Stvaranje početne populacije

Algoritam započinje nultom generacijom koja se stvara na temelju slučajnog odabira. Ova generacija nema nikakva optimizacijska svojstva, služi samo kao temelj za stvaranje i razvoj boljih jedinki u sljedećim generacijama. Postoje različiti načini kojima gradimo početnu generaciju, no najčešće korištene metode su *Full* i *Grow*.

3.1.1. Metoda full

Stvaranje stabala kod *Full* metode započinje slučajno izabranom funkcijom koja predstavlja korijen. Potom rekurzivno stvaramo sljedeće razine stabla dok ne dođemo do maksimalno predviđene dubine. Tijekom ovog procesa čvorove biramo samo iz skupa operatora. Na posljednjoj razini odabiremo čvorove isključivo iz skupa završnih znakova.

3.1.2. Metoda grow

Ova metoda je slična prijašnjoj s razlikom da u svakom trenutku stvaranja čvorove možemo birati i iz skupa završnih znakova. Gradnjom stabala na ovaj način nećemo uvijek doći do maksimalne dubine zbog čega su stabla izgrađena ovom metodom dosta raznolikija i manje složena.

3.1.3. Metoda ramped half-and-half

Najčešća metoda korištena za stvaranje nulte populacije je metoda *ramped half-and-half*. Polovica jedinki se stvara *full* metodom, a druga polovica *grow*. Određena je samo maksimalna dubina stabala, a jedinke se stvaraju za svaku dubinu, sve do maksimalne.

Prednost ove metode je što stvara jedinke s dobrom raspodjelom po veličini i strukturi.

3.2. Procjena dobrote jedinke

Procjena dobrote jedinke uvelike ovisi o problemu kojim se bavimo. Predstavljena je funkcijom dobrote (engl. *fitness function*). Ovisno o veličinama koje dobijemo primje-

nom funkcije dobrote možemo utvrditi kvalitetu pojedinog rješenja. Pravilno implementirana funkcija dobrote je važna jer o njoj u velikoj mjeri ovisi koja će rješenja ostati u razmatranju, a koja ne. Dobrotu neke jedinke možemo definirati kao mjeru kvalitete tog rješenja u zadanom prostoru rješenja. Primjerice, kod traženja maksimuma neke funkcije, dobrotu možemo prikazati kao vrijednost koju ta funkcija ima za pojedinu jedinku. Kod problema stvaranja rasporeda, dobrotu možemo procijeniti kao postotak predavanja koja se ne poklapaju, vremenskom intervalu koji studenti imaju između uzastopnih predavanja i slično.

Definiranje funkcije dobrote je jedan od ključnih problema genetskog programiranja jer je potrebno da bude što preciznija, a što jednostavnija budući da je njeno evaluiranje potrebno obavljati u svakoj generaciji na cjelokupnom skupu jedinki.

3.3. Genetski operatori

3.3.1. Selekcija

Jedan od centralnih mehanizama genetskog programiranja je selekcija. Ona omogućava boljim jedinkama da prenesu svoje "gene" u sljedeće generacije algoritma. Dobrota jedinke se zasniva na funkciji dobrote koju smo već opisali. Operatori selekcije i njihov utjecaj se mogu okarakterizirati kroz nekoliko parametara:

- vrijeme preuzimanja,
- selekcijski intenzitet,
- momenti razdiobe dobrote

Vrijeme preuzimanja je vrijeme potrebno operatoru selekcije da generira populaciju u kojoj se nalazi samo najbolje rješenje. To je vrijeme od prve generacije pa do nastanka generacije koja je popunjena samo jedinkama najboljeg rješenja ako koristimo samo mehanizam selekcije. Ako je mehanizam selekcije pravilan, u svakoj generaciji bi trebali imati sve bolja rješenja.

Selekcijski intenzitet je povećanje srednje dobrote populacije prije i poslije primjene operatora selekcije, podijeljenog sa standardnom devijacijom populacije:

$$I = \frac{\bar{f}_s - \bar{f}}{\sigma}$$

\bar{f} predstavlja prosječnu dobrotu rješenja u populaciji roditelja, \bar{f}_s predstavlja prosječnu dobrotu rješenja u populaciji djece, dok je σ standardno odstupanje u populaciji roditelja.

r -ti centralni moment razdiobe dobrote populacije rješenja veličine n se računa prema izrazu:

$$\mu_r = \frac{1}{n} \sum_{i=1}^n (f_i - \bar{f})^r$$

gdje je f_i dobrota rješenja, a \bar{f} srednja vrijednost dobrote za čitavu populaciju.

U nastavku ćemo se upoznati s nekim od često korištenih vrsta selekcije.

Proporcionalna selekcija

Proporcionalnu selekciju možemo slikovito prikazati kao kotač ruleta gdje je svakom rješenju u populaciji pridružen isječak čija je površina proporcionalna dobroti rješenja. Tako najbolja jedinka ima najveći isječak kotača i najveću šansu za odabirom, dok će kod onih manje dobrih biti obrnuto.

Turnirska selekcija

Kod turnirske selekcije iz populacije izvlačimo slučajni uzorak od s rješenja te odabiremo ono rješenje iz tog uzorka koje ima najveću dobrotu.

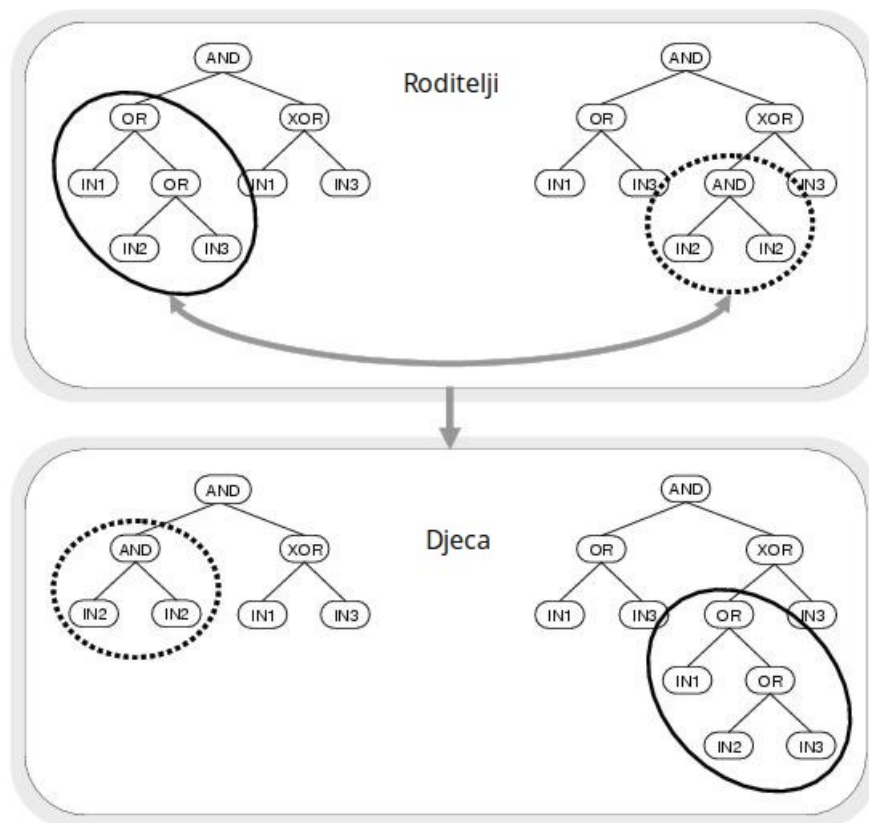
Turnirska selekcija koja u populaciji od n jedinki izvlači uzorak od njih s naziva se s -turnirska selekcija.

3.3.2. Križanje

Križanje je analogno biološkoj spolnoj reprodukciji. Križanjem dolazi do rekombiniranja genetskog materijala dvaju roditelja. Rezultat su dvije jedinke djece koje od svakog roditelja nasljeđuju dio genetskog materijala i koje će vjerojatno imati bolju dobrotu. U slučaju da se to ne dogodi loša rješenja će brzo "odumrijeti" mehanizmima selekcije koji će ih rjeđe birati za prijenos u sljedeće generacije. Ovisno o načinu prikaza jedinke postoje različite vrste križanja koje između ostaloga uključuju uniformno križanje koje se koristi kod jedinki predstavljenih kao niz bitova, aritmetičko, diskretno linearno križanje i drugi. Na slici 3.1 se nalazi primjer križanja na kojem se primjećuje kako nastala djeca nasljeđuju genetski materijal svojih roditelja.

3.3.3. Mutacija

Mutacija omogućava genetsku raznolikost jedinki čime se mogu dobiti bolja ili lošija rješenja. U generaciji se bira udio jedinki koji će sudjelovati u mutaciji ili vjerojatnost

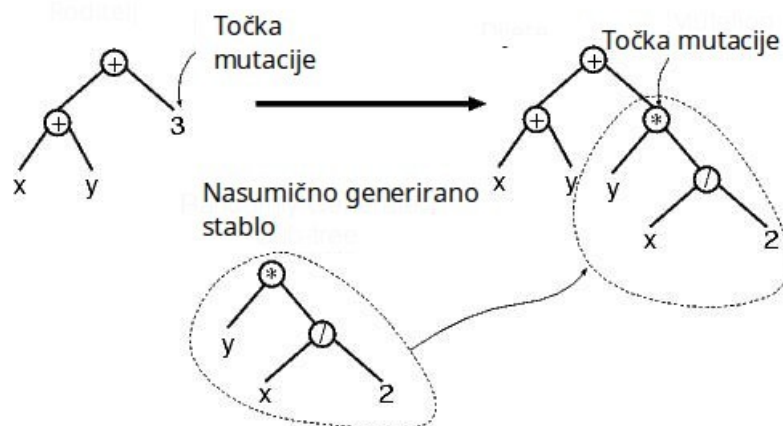


Slika 3.1: Primjer križanja

odabira mutacije za svaku pojedinu jedinku. Kod jedinki prikazanih stablima mutacija se obavlja tako da se nasumično odabere čvor koji se briše zajedno sa svojim podstablima. Na njegovom mjestu se generira novo podstablo stvoreno slučajnim odabirom. Pri tome je potrebno paziti da ne dođe do prekoračenja postavljenih ograničenja poput dubine stabla ili broja čvorova. Na slici 3.2 je prikazan primjer mutacija jedinke predstavljene stablom.

3.4. Problemi genetskog programiranja

Kao što smo već napomenuli, kod genetskog programiranja možemo naići i na određene probleme. To je prije svega problem prebrzog rasta jedinki (engl. *bloat*). Ako je jedinka prikazana u obliku stabla, problem se očituje kao nekontrolirani rast veličine stabla, to jest broja čvorova i dubine bez poboljšanja dobrote. To može dovesti do dugog trajanja evaluacije većih jedinki i nerazumljivosti i nečitljivosti rješenja. Predloženi su različiti mehanizmi rješavanja tog problema koji između ostalog uključuju kažnjavanje prevelikih jedinki, odbacivanje onih jedinki koji prelaze dopuštena ogra-



Slika 3.2: Mutacija jedinke predstavljene stablom

ničenja te njihova zamjena novim ili podrezivanje stabala (engl. *tree pruning*).

3.5. Primjena genetskog programiranja na problem klasifikacije

U ovom potpoglavlju ćemo se detaljnije upoznati s načinom na koji se genetsko programiranje koristi za rješavanje problema klasifikacije podataka. Klasifikacija je raširen problem s kojim se i nesvjesno suočavamo svaki dan. Na primjer, prepoznavanje lica pojedinaca koje smo prije upoznali, vrsta životinja, biljaka i hrane i drugih objekata iz okoline. Koristi se i u složenim znanstvenim domenama poput medicine gdje može pomoći u dijagnozi pacijenata na osnovu povijesti bolesti drugih pacijenata sa sličnim simptomima.

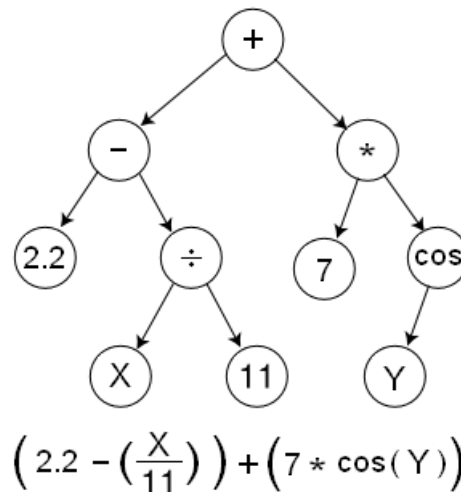
S obzirom na raširenost ovog problema važno je razviti računalne mehanizme sposobne rješavanja problema klasifikacije u slučajevima kada je to za ljude presloženo i vremenski zahtjevno. Generalni oblik klasificiranja podrazumijeva prepoznavanje klase kojima određeni primjerci problema pripadaju na osnovu atributa ili svojstava ispitnih primjeraka. Na osnovu ispitnih primjeraka možemo razviti model koji će predviđati klasu kojoj pripadaju novi, još nesvrstani primjerci.

Kroz povijest razvoja klasifikacije koristili su se različiti algoritmi i načini rješavanja tog problema. Genetsko programiranje je relativno novo i brzo razvijajuće područje koje se u široj primjeni našlo tek u zadnjih nekoliko desetljeća.

Model na kojem ćemo se zadržati u ovom radu je reprezentacija regresijskim stablima koje koristimo za klasifikaciju numeričkim izrazima. U kombinaciji s aritme-

tičkim operacijama koristit ćemo se i logičke izraze poput *if-then-else* operacija koje se često koriste u računalnim programima. Ovakav način prikaza rješenja se pokazao relativno uspješnim u praktičnim primjenama poput klasifikacije rukopisa, bolesti, prepoznavanja objekata i slično.

Primjer stabla i pripadajuće funkcije koja ovisi o ulaznim parametrima X i Y prikazan je na slici 3.3



Slika 3.3: Stablo odluke

Klasifikator numeričkih izraza se definira kao bilo koja funkcija koja vraća numeričku vrijednost. Funkcija se sastoji od matematičkih operacija, uključujući aritmetičke operacije i druge poput sinusa, kosinusa ili logaritma. Kao što smo prije napomenuli, funkcije mogu sadržavati i operatore usporedbe kao ako-onda-ili izraz. Klasifikator numeričkih izraza provodi klasifikaciju pomoću varijabli koje su mu predane kao ulaz i na izlazu daje numeričku vrijednost. Ta numerička vrijednost se onda mora interpretirati u jednu od mogućih klasa za zadani problem.

Važnost će biti na preciznosti modela koji nastaju genetskim programiranjem. Iako postoje drugi faktori o kojima može ovisiti kvaliteta promatranog rješenja poput veličine programa, razumljivosti, vremenskoj izvedbi i slično, njih ćemo svrstati u drugi plan. Na osnovu toga, kvalitetna rješenja će biti ona koja imaju najveći preciznost izvođenja klasifikacije, bez obzira na ostale promatrane uvjete.

Cilj je za zadane primjere proizvesti program koji će što preciznije klasificirati primjere i u velikom postotku predviđati pripadnost primjeraka problema zadanim klasama.

4. Programsko ostvarenje

U ovom poglavlju ćemo detaljnije opisati programsko ostvarenje i opisati na koji smo način rješavali problem klasifikacije podataka upotrebom genetskog programiranja. Upoznati ćemo se s radnim okruženjem ECF-a, opisati genetske operatore koje smo koristili i objasniti prikaz dobivenog rješenja.

4.1. Radno okruženje ECF

ECF (engl. *Evolutionary Computation Framework*) je radno okruženje za rješavanje problema evolucijskog računanja koje se razvija na Fakultetu elektrotehnike i računarstva u Zagrebu pod vodstvom izv. prof. dr. sc. Domagoja Jakobovića. Napisano je u programskom jeziku C++.

Za korištenje ECF-a možemo kombinirati i parametrizirati sljedeće komponente:

- Genotype
- Algorithm
- Evolutionary system

Genotype je osnovna komponenta ECF-a. Svaka jedinka u populaciji ima jednog ili više genotipa koje je najjednostavnije definirati u konfiguracijskoj datoteci. Genotip je jedina komponenta koju je potrebno definirati, sve ostale mogu koristiti pretpostavljene vrijednosti. Svaki genotip ima svoju strukturu i genetske operatore križanja i mutacije. U ovom radu ćemo koristiti izvedeni razred **Tree** koji predstavlja stablasti kromosom genetskog programiranja.

Algorithm predstavlja algoritam koji koristimo. Željeni algoritam i njegovi parametri se mogu definirati u konfiguracijskoj datoteci.

Evolutionary system je kontekst u kojem postoje genotip i algoritam. Kontekst se brine za sve ostalo: veličinu populacije, uvjete završavanja i slično.

Parametri se zadaju u konfiguracijskoj datoteci XML formata (engl. *Extensible Markup Language*), a dohvaćaju se prilikom inicijaliziranja ECF-a. Primjer konfigu-

racijske datoteke dan je na slici 4.1

```
<ECF>
  <Algorithm>
    <SteadyStateTournament>
      <Entry key="tsize">3</Entry>
    </SteadyStateTournament>
  </Algorithm>
  <Genotype>
    <FloatingPoint>
      <Entry key="lbound">-50</Entry>
      <Entry key="ubound">50</Entry>
      <Entry key="dimension">3</Entry>
    </FloatingPoint>
  </Genotype>
  <Registry>
    <Entry key="function">6</Entry>
    <Entry key="population.size">500</Entry>
    <Entry key="mutation.indprob">0.3</Entry>
    <Entry key="term.maxgen">500</Entry>
    <Entry key="term.fitnessval">1e-12</Entry>
    <Entry key="log.frequency">10</Entry>
    <Entry key="log.level">3</Entry>
  </Registry>
</ECF>
```

Slika 4.1: Primjer konfiguracijske datoteke

4.2. Korišteni klasifikatori

U ovom radu ćemo implementirati dva različita klasifikatora genetskog programiranja. To su :

- Regresijski GP (klasifikacijski aritmetički izrazi)
- Regresijski GP uz glasanje - posebno stablo za svaku klasu

Regresijski GP s jednim stablom uspoređuje numeričke vrijednosti dobivene evaluacijom stabla s intervalima definiranim za svaku klasu. Razred kojem primjer pripada je onaj u čijem se intervalu nalazi dobivena numerička vrijednost. Funkcijom dobrote potom određujemo uspješnost predviđanja klasa svake jedinice.

Kod regresijskog GP-a uz glasanje svaka klasa ima definirano svoje stablo. Tijekom evaluacije računamo izlaznu numeričku vrijednost za svako od stabala. Odlučujuće je ono stablo koje daje najmanju apsolutnu vrijednost.

Tada u slučaju da imamo tri klase, dakle tri stabla, možemo dobiti sljedeće vrijednost za svako od stabala:

stablo A: 3

stablo B: 1.5

stablo C: -2

Recimo da primjer odgovara klasi B. Odlučujuće stablo je stablo B budući da upravo ono ima najmanju apsolutnu vrijednost. Time je jedinka dobro klasificirala zadani primjer.

4.3. Funkcije čvorova

Stabla odluke se sastoje od skupa funkcijskih znakova i skupa završnih znakova. U skup završnih znakova spadaju atributi klasifikacijskog problema. Ako je ukupan broj atributa jednak n varijable tada imaju oznake od x_0 do x_{n-1} . Prilikom evaluacije, za vrijednosti varijabli iz skupa završnih znakova uzimaju se odgovarajuće vrijednosti iz skupa podataka definiranih za svaki problem klasifikacije.

U skup funkcijskih znakova spadaju:

- Aritmetički operatori
- Sinus, kosinus
- Max, min
- Sqrt, Log
- IfPositive
- IfLessThanEq

U nastavku ćemo se pobliže upoznati s njima.

U aritmetičke operatore koje koristimo spadaju zbrajanje (+), oduzimanje (-), množenje (*) i dijeljenje (/). Pored njih, koristit ćemo i funkcije sinusa i kosinusa.

Ostale korištene funkcije je potrebno implementirati unutar radnog okruženja ECF-a. Pri tome im je potrebno definirati pridružena svojstva **nArguments_** koje označava broj djece koje će čvor s navedenom funkcijom imati te **name_** koje označava ime navedene funkcije. U svakoj definiranoj funkciji je potrebno implementirati metodu **execute** u kojoj vrijednosti čvorova djece dohvaćamo pomoću metode **getNextArgument**. Izlazna vrijednost čvora je vrijednost dobivena primjenom odgovarajuće funkcije na jedno ili više čvorova djece.

Max i *Min* funkcije definiraju čvorove koji imaju dvoje djece. *Max* funkcija vraća vrijednost onog djeteta koje ima veću vrijednost, dok *Min* funkcija vraća manju od vrijednosti svoje djece.

Funkcija *Sqrt* definira čvor koji ima jedno dijete te vraća korijen apsolutne vrijednosti tog djeteta.

Kod funkcije *Log* je potrebno pripaziti na vrijednost koju joj predajemo. Budući da je logaritamska funkcija definirana samo za pozitivne vrijednosti veće od nule, u slučaju da je vrijednost djeteta manja ili jednaka nuli funkcija vraća logaritam od jedan, to jest nula.

IfPositive definira čvor koji ima troje djece. U slučaju da je vrijednost prvog djeteta veća od 0 vraća vrijednost drugog djeteta, a inače vraća vrijednost trećeg.

IfLessThanEq definira čvor koji ima četvero djece. U slučaju da je vrijednost prvog djeteta manja ili jednaka vrijednosti drugog vraća vrijednost trećeg djeteta, a u suprotnom slučaju vraća vrijednost četvrtog djeteta.

Uvođenje logičkih izraza pored onih aritmetičkih je važno jer omogućava veću raznolikost rješenja i pomaže u preciznijoj evoluciji stabala.

4.4. Funkcije dobrote

Funkcija dobrote koju ćemo koristiti u treniranju klasifikatora je *F1 score*. Nastala je iz matrice (*confusion matrix*) koja omogućava vizualizaciju uspješnosti algoritma, obično nadziranog učenja kod problema umjetne inteligencije među koje spada i problem klasifikacije.

Sastoji se od dvije dimenzije ("stvarno" i "predviđeno") i identičnog skupa klasa u obje dimenzije.

Za primjer možemo uzeti klasifikator koji razvrstava pse, mačke i zečeve u njihove pripadajuće klase. Ako uzmemo skup od 27 životinja - 8 mačaka, 6 pasa i 13 zečeva, rezultati nakon klasifikacije bi mogli izgledati ovako:

		Predviđeno		
		Mačka	Pas	Zec
Stvarno	Mačka	5	3	0
	Pas	2	3	1
	Zec	0	2	11

Od osam mačaka, klasifikator je za tri predvidio da su psi, a od šest pasa za jednog je predvidio da je zec, a za dva da su mačke. Svi pogoci su prikazani dijagonalom pa je lako vidjeti sve pogreške.

Matrica iz koje ćemo dobiti funkciju dobrote ima dva retka i dva stupca sa sljedećim mogućim vrijednostima za klasu koju analiziramo:

- *false positive (FP)* - Svrstavanje pogrešnog primjerka u klasu

- *false negative (FN)* - Primjerci klase koji su svrstani u neku drugu klasu
- *true positive (TP)* - Pogodak
- *true negative (TN)* - Točno neprihvatanje primjerka u klasu

Za navedeni primjer i klasu Mačka, dobili bi sljedeću matricu.

5 TP (točno svrstane mačke)	3 FN (mačke svrstane u pse)
2 FP (psi svrstani u mačke)	17 TN (sve ne-mačke svrstane u druge klase)

Naša funkcija dobrote će tada imati oblik:

$$F1 = \frac{2 * TP}{2 * TP + FP + FN}$$

Moguće vrijednosti su iz intervala $[0, 1]$ gdje veći broj predstavlja veću dobrotu rješenja. Za probleme s više klasa uzima se prosječna vrijednost funkcije dobrote za sve klase.

4.5. Parametri i konfiguracijske datoteke

Parametre zadajemo u konfiguracijskoj datoteci XML formata koju je potrebno predati prilikom inicijaliziranja ECF-a. Tu možemo definirati mehanizme križanja, veličinu populacije, vjerojatnost mutacije, minimalnu i maksimalnu dubinu stabla i druge.

U slučaju ne navođenja nekog od parametara, uvode se sljedeće pretpostavljene vrijednosti:

- Veličina populacije: 100
- Vjerojatnost mutacije: 0.3
- Prekid: 50 generacija bez napretka najbolje jedinke

Parametri genotipa koji će nama biti od značaja su *functionset* pomoću kojeg definiramo skup funkcijskih znakova te *terminalset* koji predstavlja skup završnih znakova. Skup završnih znakova ćemo u svim slučajevima definirati oznakama od x_0 do x_{n-1} gdje n predstavlja broj atributa problema.

Pored navedenih možemo definirati i parametar *classesNum* koji predstavlja broj klasa koji koristimo u klasificiranju. U slučaju ne navođenja navedenog parametra koristi se pretpostavljena vrijednost dva.

Primjer parametarske datoteke gdje možemo vidjeti navedene atribute nalazi se na slici 4.2.

```

<ECF>
  <Genotype>
    <Tree>
      <Entry key="maxdepth">4</Entry>
      <Entry key="mindepth">1</Entry>
      <Entry key="functionset">sqrt log iflte max min ifpos sin cos + - * </Entry>
      <Entry key="terminalset">X0 X1 X2 X3</Entry>
    </Tree>
  </Genotype>
  <Registry>
    <Entry key="classesNum">2</Entry>
  </Registry>
</ECF>

```

Slika 4.2: Primjer konfiguracijske datoteke

Kod definiranja intervala klasa koristimo datoteku definiranu parametrom *classesfile*. U slučaju ne navođenja parametra koristi se pretpostavljeni naziv *classes.txt*. Pomoću ove datoteke definiramo sve klase koje koristimo kao i intervale kojima pojedina klasa pripada. Ako datoteka ne postoji, koriste se pretpostavljene klase s oznakama od 0 do $n - 1$ gdje n predstavlja broj klasa definiranih parametrom *classesNum*. Za intervale se koriste negativna i pozitivna potencija broja dva, gdje eksponent intervala odgovara numeričkoj oznaci klase. U klasi $n - 1$ definiramo dva intervala koja pokrivaju sve preostale realne vrijednosti.

Naprimjer, ako smo definirali postojanje tri klase bez datoteke sa konfiguracijom dobit ćemo sljedeće vrijednosti intervala:

- Klasa 0: $[-1, 1]$
- Klasa 1: $[-2, -1] \cup [1, 2]$
- Klasa 2: $(-\infty, -2] \cup [2, +\infty)$

Primjer strukture datoteke za definiranje razreda nalazi se na slici 4.3.

```

a 1:4 6:12
b :1 4.5:6
c 4:4.5 12:

```

Slika 4.3: Primjer definiranja klasa

Na početku retka definirano je ime klase, a intervale odvajamo dvotočkama. Beskonačne vrijednosti označavamo ne navođenjem vrijednosti s odgovarajuće strane intervala.

Za primjer naveden na slici definirane su tri klase i intervali:

- Klasa a: $[1, 4] \cup [6, 12]$
- Klasa b: $(-\infty, 1] \cup [4.5, 6]$
- Klasa c: $[4, 4.5] \cup [12, +\infty)$

Primjere za učenje definiramo u datoteci definiranoj parametrom *inputfile* s pretpostavljenom vrijednosti *learning.txt* u kojoj navodimo attribute svakog primjera i razred kojem pripada.

Svaki primjer je napisan u novom retku, a atributi su odvojeni zarezima. Kod primjera s n atributa imat ćemo $n + 1$ zapisa u retku gdje posljednji predstavlja klasu.

Na slici 4.4 prikazana je datoteka za učenje čiji primjeri imaju četiri atributa i dva razreda.

```
2,1,0,0,1
0,0,1,1,0
2,2,0,0,1
1,2,0,1,1
0,1,1,0,0
1,0,0,0,1
0,2,0,0,1
1,1,1,1,1
2,2,0,1,0
2,1,1,0,1
0,1,0,1,1
1,0,1,0,1
2,1,1,1,0
0,0,1,0,0
```

Slika 4.4: Datoteka sa primjerima za učenje

Datoteka *test.txt* je pretpostavljane vrijednost parametra *testfile* u kojoj se nalaze primjeri za testiranje koji imaju identičnu strukturu.

Detaljnije objašnjenje o atributima i šta točno znače za navedene datoteke ćemo opisati u sljedećem poglavlju.

Ime datoteke u koju se zapisuju rezultati određeno je parametrom *resultsfile* koji ima pretpostavljenu vrijednost *results.txt*.

U datoteci se za svaku generaciju zapisuje broj generacije i najveće vrijednosti dobrote skupa za učenje i testiranje. Vrijednosti su odvojene zarezima. Primjer je prikazan na slici 4.5.


```
Gen_No, Training, Test
0,0.837743,0.888288
1,0.837743,0.888288
2,0.837743,0.888288
3,0.837743,0.888288
4,0.837743,0.888288
5,0.837743,0.888288
6,0.882051,0.935214
7,0.882051,0.935214
8,0.882051,0.935214
9,0.903199,0.956944
10,0.932367,0.978495
11,0.941919,1
12,0.941919,1
13,0.942017,1
14,0.942017,1
15,0.951528,0.978405
```

Slika 4.5: Struktura datoteke sa rezultatima

5. Primjena i analiza rezultata

Implementirane klasifikatore ćemo primijeniti na dva primjera problema klasifikacije. Prvi je vjerojatno najpoznatiji skup podataka za problem klasifikacije u kojem je cilj klasificirati različite porodice cvijeta iris. Drugi problem je problem klasifikacije stakla na osnovu svojstava i kemijskog sastava.

5.1. *Iris flower data set*

Iris flower data set je vjerojatno najpoznatiji skup podataka koji se koristi kod problema klasifikacije. Predstavio ga je Ronald Fisher 1936. godine u svom radu *The use of multiple measurements in taxonomic problems*.

Podaci se sastoje od tri klase po 50 primjeraka gdje je svaki primjerak definiran sa četiri atributa. Klase predstavljaju tri porodice cvijeta iris - porodice su *iris setosa*, *iris versicolor*, *iris virginica*.

Atributi koji se koriste su:

- dužina lapa
- širina lapa
- dužina latica
- širina latica

Sve vrijednosti su predstavljene u centimetrima.

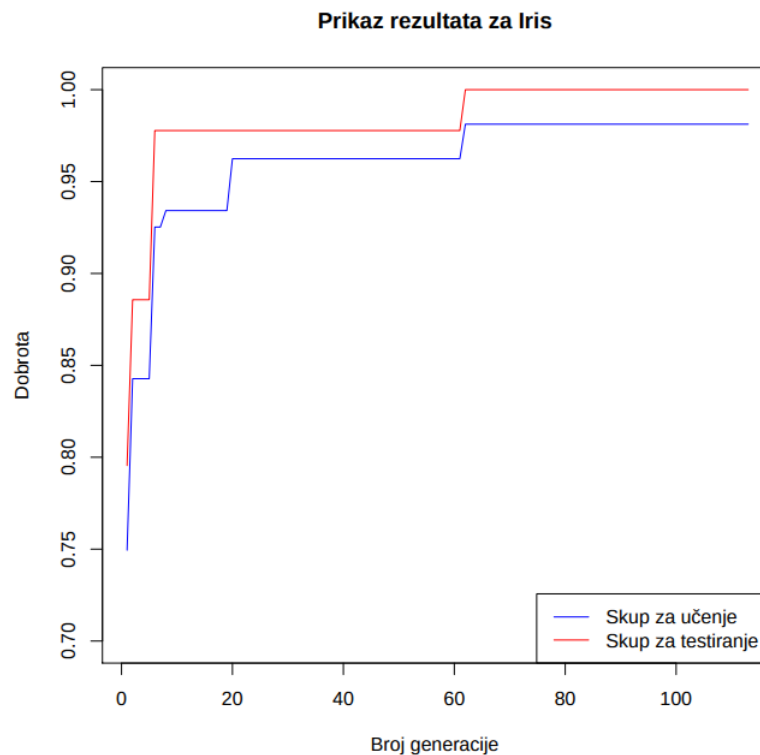
Nekoliko uzoraka je prikazano sljedećom tablicom.

Dužina lapa	Širina lapa	Dužina latica	Širina latica	Porodica	Oznaka
5.1	3.5	1.4	0.2	<i>I. setoosa</i>	0
4.9	3.0	1.4	0.2	<i>I. setosa</i>	0
7.0	3.2	4.7	1.4	<i>I. versicolor</i>	1
6.4	3.2	4.5	1.5	<i>I. versicolor</i>	1
6.3	3.3	6.0	2.5	<i>I. virginica</i>	2
5.8	2.7	5.1	1.9	<i>I. virginica</i>	2

Primjer ćemo testirati na populaciji od 100 jedinki, s dubinom stabla između 1 i 5 te faktorom mutiranja 0.3. U svim rezultatima omjer primjera za učenje i ispitivanje je 70:30.

Tijekom učenja klasifikatora ujedno ćemo u svakoj generaciji ispitati i uspješnost klasificiranja na skupu za ispitivanje.

Rezultate primjenom regresijskog GP-a s jednim stablom možemo vidjeti na slici 5.1



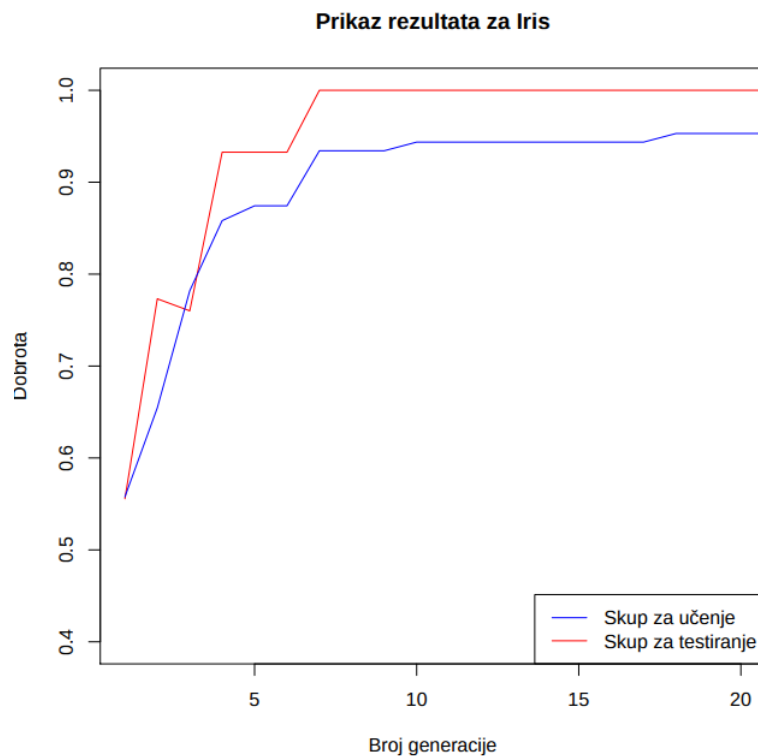
Slika 5.1: Klasifikator s jednim stablom

Možemo zaključiti da klasifikator jako brzo dolazi do zadovoljavajućeg rješenja pomoću kojeg može sa 100% uspješnosti predvidjeti klasu kojoj pripadaju primjeri za testiranje.

Kod regresijskog GP-a s više stabala vidimo da postizemo slične rezultate no da klasifikator uz dosta manje generacija dolazi do zadovoljavajućih rezultata.

Budući da se radi o problemu sa samo tri klase i početna populacija postiže rezultate s kojima može klasificirati primjere za testiranje s uspješnosti većom od 50%.

Kako bi detaljnije testirali uspješnost rada klasifikator, primjeniti ćemo ih i na problem u kojem je cilj raspodijeliti primjere između većeg broja klasa.



Slika 5.2: Klasifikator s više stabala

5.1.1. Klasifikacija stakla

Sljedeći problem na kojem ćemo testirati klasifikatore je problem identifikacije vrste stakla. Američka agencija *Forensic Science Service* je testirala svojstva i kemijski sastav sedam vrsta stakala.

U svojoj analizi testirali su sljedeća svojstva i udjele sastojaka:

- Lom svjetla
- Natrij, Magnezij, Aluminijski
- Silicij, Kalij, Kalcij
- Barij, Željezo

Navedeno će ujedno biti atributi kojima ćemo učiti klasifikatore.

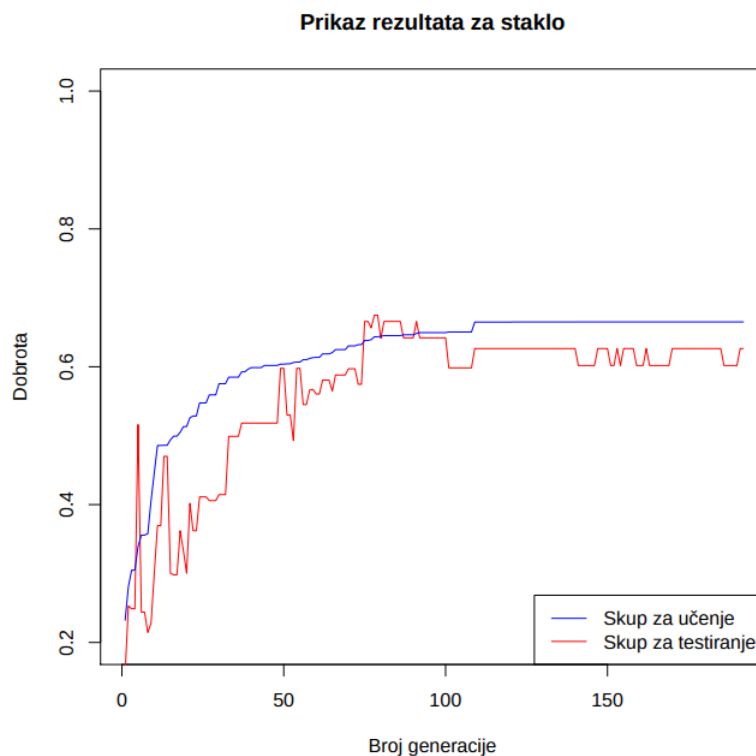
Podaci se sastoje od 6 vrste stakla od kojih svi nemaju jednak broj primjera:

- Obrađeno staklo prozora zgrada - 70 primjeraka
- Neobrađeno staklo prozora zgrada - 75 primjeraka
- Obrađeno staklo prozora vozila - 17 primjeraka
- Staklo posuda - 13 primjeraka

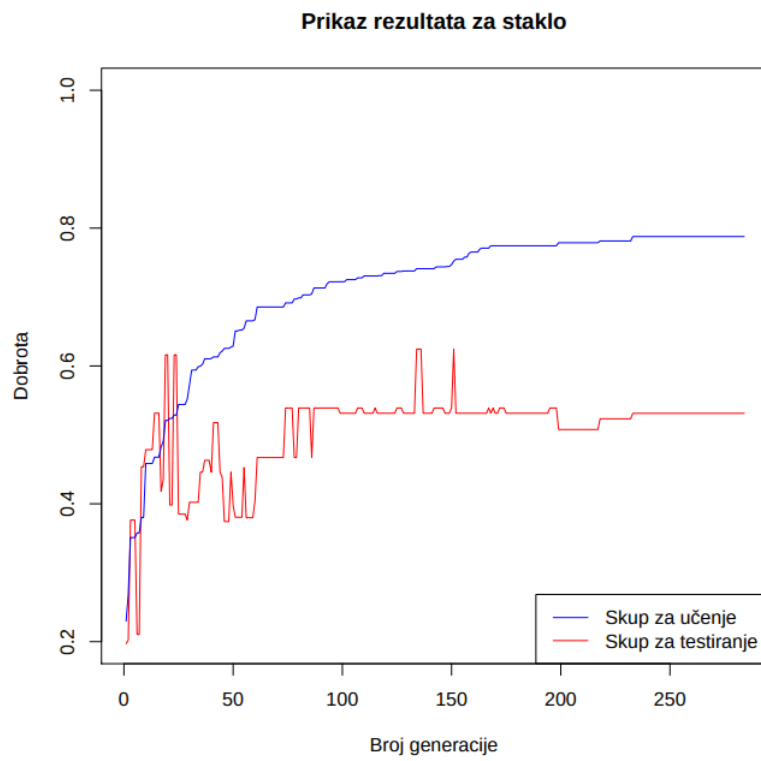
- Stolno posuđe - 9 primjeraka
- Staklo automobilskih farova - 29 primjeraka

Vidimo da se problem sastoji od većeg broja atributa i klasa, te da su brojevi primjeraka za pojedinu klasu dosta neuravnoteženi zbog čega možemo očekivati da će klasifikatorima biti puno teže pravilno odrediti klase.

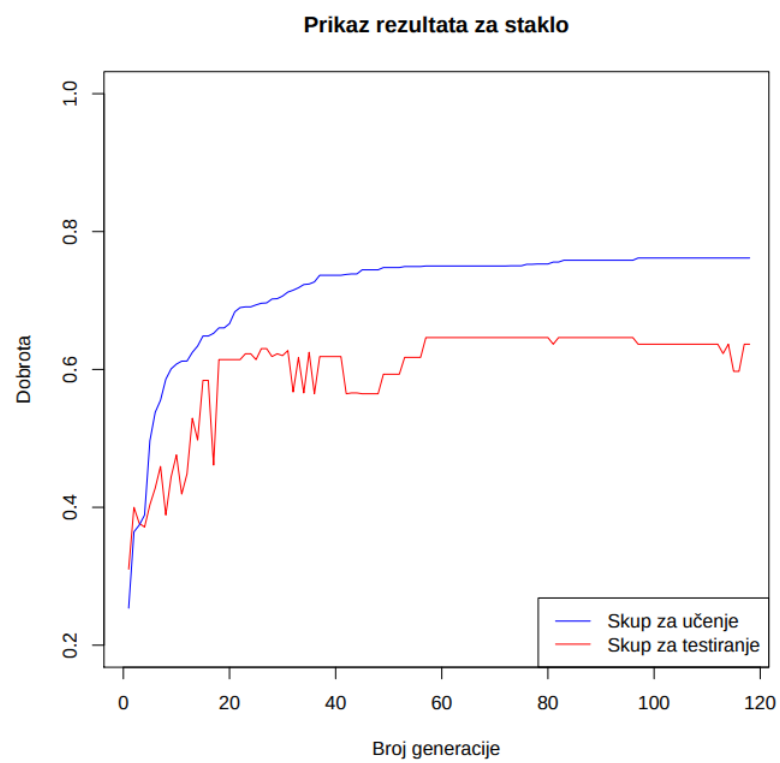
Dubina stabla će ponovno biti između 1 i 5, no za ovaj problem ćemo koristiti populaciju od 500 jedinki. Problem ćemo testirati i s različitim vjerojatnostima mutacije.



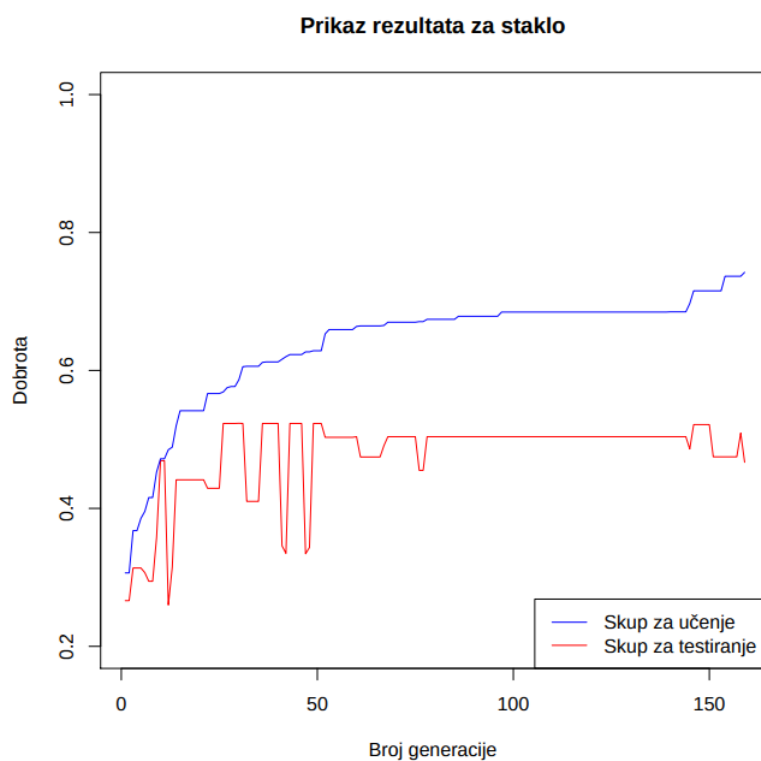
Slika 5.3: Klasifikator s jednim stablom i faktorom mutacije 0.3



Slika 5.4: Klasifikator s jednim stablom i faktorom mutacije 0.9



Slika 5.5: Klasifikator s više stabala i faktorom mutacije 0.3



Slika 5.6: Klasifikator s više stabala i faktorom mutacije 0.9

Ponovno možemo primijetiti da klasifikator s više stabala dolazi do sličnih ili boljih rezultata u značajno manjem broju generacija. Veći faktor mutacije obično dovodi i do boljih rezultata skupa za testiranje, no vidimo da to obično ne pomaže točnosti klasifikacije primjera za učenje. U ovom problemu klasifikatoru veliki problem predstavlja neuravnotežen broj primjera za pojedine klase zbog čega ne uspijevamo dobiti veliku točnost. Na rezultate može utjecati i izbor primjera koje koristimo za testiranje i učenje zbog čega bi bila potrebna detaljna analiza utjecaja svojstava na rezultate.

6. Zaključak

U radu smo opisali genetsko programiranje i problem klasifikacije te se malo поближе upoznali sa načinom na koji rade. Vidjeli smo da se genetsko programiranje može uspješno koristiti za učenje klasifikatora no da je potrebno pripaziti na moguće probleme poput pojavljivanja pretjerane složenosti i prevelikog broja čvorova u stablima. Problemi s više atributa i klasa zahtijevaju i veće računalne resurse zbog čega ponekad nisu najbolje rješenje za velike izračune bez pristupa jačim uređajima. Navedeni klasifikatori se mogu dodatno poboljšati primjenom različitih operatora križanja mutiranja kao i funkcija koje se koriste u čvorovima za izračune. Moguće je i uvođenje skupa podataka kojim bi osigurali da ne dođe do prilagođavanja klasifikatora specifičnim podacima i šumu.

LITERATURA

1. Riccardo Poli, William B. Langdon, Nicholas F. McPhee, A Field Guide to Genetic Programming.
http://www.cs.ucl.ac.uk/staff/wlangdon/ftp/papers/poli08_fieldguide.pdf
2. Mehryar Mohri, Afshin Rostamizadeh, Ameet Talwalkar, Foundations of Machine Learning.
<http://www.cs.nyu.edu/mohri/mlbook/>
3. Wei-Yin Loh, Classification and regression trees.
<http://www.stat.wisc.edu/loh/treeprogs/guide/wires11.pdf>

Primjena genetskog programiranja na problem klasifikacije podataka

Sažetak

Problemi klasifikacije su danas jedni od najraširenijih problema s kojima se suoćemo. U radu se opisuje uporaba genetskog programiranja na problem klasifikacije te implementacija i primjena klasifikatora s jednim ili više stabala odluke. Napravljena je analiza njihovog rada te su opisani problemi i moguća poboljšanja.

Ključne riječi: ECF, genetsko programiranje, klasifikacija, stabla odluke

Classification of data with genetic programming

Abstract

Classification problems are one of the widely spread problems that are encountered in daily life. Our goal in this thesis is to describe use of genetic programming on classification and implementation of the classifiers with one or more decision trees together with detailed analysis and explanation of possible problems and improvements.

Keywords: ECF, genetic programming, classification, decision trees