

SVEUČILIŠTE U ZAGREBU
FAKULTET ELEKTROTEHNIKE I RAČUNARSTVA

DIPLOMSKI RAD br. 1479

**SUSTAV ZA DETEKCIJU PLAGIJATA TE
ODREĐIVANJE AUTORSTVA IZVORNIH
KODOVA**

Dino Rakipović

Zagreb, lipanj 2017.

Sadržaj

Uvod	2
1 Određivanje autorstva	4
1.1 Izvlačenje značajki	6
1.1.1 Leksičke značajke	6
1.1.2 Strukturne značajke	7
1.1.3 Sintaksne značajke	8
1.2 Selekcija značajki	8
1.2.1 Selekcija značajki po sadržaju informacije	9
1.2.2 Selekcija značajki po iznosu varijance	9
1.3 Slučajna šuma	9
1.3.1 Gini nečistoća	10
1.3.2 Uzajamni sadržaj informacije	10
1.4 Prikupljanje podataka	11
1.5 Rezultati i rasprava	11
1.5.1 Rezultati bez selekcije značajki	11
1.5.2 Rezultati uz selekciju značajki sadržajem informacije	12
1.5.3 Rezultati uz selekciju značajki varijancom	13
1.5.4 Rasprava	13
1.6 Implementacijski detalji	14
2 Određivanje sličnosti izvornih kodova	15
2.1 Izračun otiska izvornog koda	16
2.1.1 Rabin-Karp algoritam	16
2.1.2 Winnowing	17
2.2 Preoblikovanje izvornog koda	18

2.2.1	Zamjena imena varijabli	18
2.2.2	Zamjena while petlji u for petlje	19
2.3	Izračun sličnosti izvornih kodova	19
2.4	Usporedba sa starim sustavom	20
3	Integracija dva sustava	21
3.1	Korisničko sučelje	21
3.2	Implementacijski detalji	24
	Zaključak	25
	Literatura	26
A	Tipovi čvorova apstraktnog sintaksnog stabla	29
B	Primjer izvornog koda prije i poslije automata	30

Uvod

Plagijat(eng. *Plagiarism*) ili postupak krađe tuđeg rada postaje sve veći problem u današnjem svijetu te ga pronalazimo akademskom(npr. eseji, znanstveni članci) i neakademskom(npr. knjige, pjesme) svijetu. Razlog tomu je što količina podataka na internetu, koji je postao dostupan velikom broju ljudi, raste eksponencijalno te je vrlo lagano ukrasti tuđi rad i predstaviti ga kao vlastiti. U ovom radu naglasak će biti na detekciji plagijata izvornih kodova. Cilj je izraditi sustav koji bi ubrzao i uvelike pomogao u detekciji plagijata među izvornim kodovima ljudima koji se brinu o programerskim natjecanjima, laboratorijskim vježbama itd.

Detekciji pristupamo iz dva kuta, jedan je određivanje autora izvornog koda tzv. deanonimizacija autorstva, a drugi određivanje sličnosti među parovima izvornih kodova. Određivanje sličnosti parova izvornih kodova je već opisano i implementirano u završnom radu autora ovog diplomskog rada, no ovdje predstavljam brže te arhitekturno ljepše rješenje istog problema. Za određivanje autorstva bitno je primjetiti da svaki autor ostavlja svoj jedinstveni otisak dok piše programski kod, barem se tome nadamo. Kako bi autore razlikovali korištene su tehnike strojnog učenja u konkretnom slučaju klasifikator nasumične šume. Korisnik mora najprije predati skup za treniranje, koji se sastoji od izvornih kodova te točno označenih autora, klasifikatoru kako bi kasnije mogao utvrditi autorstvo na skupu koji želi provjeriti. Korisnika se prepoznaje po njegovom stilu programiranja tako da se izvorni kod pretvori u vektor brojeva od kojih je svaki broj nekakva unaprijed označena značajka(npr. ostavlja li autor novi red prije otvaranja vitičastih zagrada, koliko naredbi grananja koristi, itd.).

Dva pristupa su na kraju integrirana pod istim web sučeljem nazvanim *Turtle*. *Turtle* je prva, prema autorovom najboljem znanju, web aplikacija koja objedinjuje ova dva pristupa. Ovo sučelje nudi uvid u sve parove izvornih kodova za koje sumnja da su plagijati te ispisuje njihove sličnosti. Također kako bi detekcija bila još lakša i brža nudi uvid u izvorne kodove svakog para te boja slične dijelove jednakim bojama. Također ispisuje autore za koje misli da imaju najveću vjerojatnost da su napisali promatrani kod. Utvrđivanje autorstva je veliki korak naprijed nad rješenjem napisanim za završni rad jer nam omogućuje detekciju plagijata među raznim akademskim godinama ukoliko se laboratorijski zadaci mijenjaju. Za kraj treba naglasiti da sustav trenutno podržava samo jezik C++, ali je u implementacija ostavljena mogućnost lakog dodavanja novih jezika.

POGLAVLJE 1

Određivanje autorstva

U svijetu u kojem ne postoje unaprijed određena pravila pisanja programskog koda možemo pretpostaviti da svaki programer ostavlja svoj jedinstveni otisak dok programira. Cilj nam je kreirati klasifikator koji bi nam mogao odvojiti autore prema njihovom stilu programiranja. Ovakav klasifikator bi bilo moguće primjeniti na raznim open source projektima na kojima autori razvijaju kod anonimno te bi takav klasifikator mogao narušiti privatnost programera, ali ipak u ovom radu veći naglasak je na detekciji plagijata izvornih kodova te ovakav klasifikator koristimo nad laboratorijskim vježbama na fakultetima ili na nekim programerskim natjecanjima gdje autori predaju kod pod svojim imenom.

Za rješavanje ovog problema korišteno je strojno učenje. Strojno učenje je grana umjetne inteligencije koja se bavi algoritmima koji mogu učiti na i raditi predviđanja nad skupovima podataka. Kako bi mogli strojno učiti moramo imati skupove podataka s označenim kategorijama nad kojima algoritam uči. U ovom slučaju podaci su izvorni kodovi, a kategorije autori koji su ih napisali. Konkretno, korišten je klasifikator slučajne šume. Konfiguracija klasifikatora je detaljnije opisana u nastavku poglavlja. Klasifikacija je postupak u kojem određujemo kojoj kategoriji (od unaprijed određenih) novi podaci pripadaju. Algoritmi strojnog učenja uglavnom primaju ulazne podatke u obliku brojeva pa je potrebno izvorni kod pretvoriti u vektor brojeva u kojem će svaki broj biti neka od značajki. Te značajke su podijeljene u leksičke, sintaksne i struktu-

ralne.

Što bolje značajke odaberemo algoritam će bolje moći odvajati kategorije tj. autore. Značajke dobijemo parsiranjem izvornog koda te ću postupak detaljnije opisati u nastavku poglavlja.

1.1 Izvlačenje značajki

Kao što je već spomenuto, kako bi algoritmi strojnog učenja radili potrebni su im brojčani podaci kao ulazi. Izvorni kod se u brojčani vektor značajki pretvara koristeći ideju prvi put opisanu u radu [1], a ideja je da se izvorni kod pretvori u vektor značajki sastavljen od tri dijela, leksičkog, sintaksnog i strukturnog. Leksičke i strukturalne značajke se dobiju izravno parsiranjem izvornog koda, dok nam je za sintaksne značajke potrebno apstraktno sintakšno stablo izvornog koda. Ovako definiran skup značajki je drugačiji za svaki pojedini programski jezik zbog različitosti među njima (npr. drugačije ključne riječi) te je potrebno napisati poseban parser za svaki od njih. U ovom radu naglasak je na programskom jeziku C++ te je izvlačenje značajki implementirano samo za njega.

U nastavku su detaljno opisana i objašnjena sva tri tipa značajki. Većina tih značajki preuzeta je iz [1] dok su neke ideja samog autora. U većini značajki korištena je matematička operacija prirodnog logaritmiranja zbog svojstva da kako idemo prema većim vrijednostima ona sve manje i manje raste te dobro opisuje relativne razlike među značajkama te su neke podijeljene s duljinom izvornog koda kako bi bolje opisale frekvenciju.

1.1.1 Leksičke značajke

Leksičke značajke opisuju preferira li autor izvornog koda neke ključne riječi više od drugih (npr. `for` više od `while`), koristi li više funkcije ili piše monolitan kod, razne statistike (npr. prosječan broj parametara unutar funkcija), itd. Također izvorni kod se tokenizira te se računa frekvencija tako dobivenih tokena. *Tablica 1.1* detaljno opisuje svaku od korištenih značajki.

Tablica 1.1: Definicija leksičkih značajki

<i>Ime značajke</i>	<i>Definicija</i>	<i>Izraz</i>	<i>Veličina</i>
Frekvencija unigrama	Unigram definiramo kao jednu riječ izvornog koda	UnigramFreq	dinamično, oko 20000 za 2160 izvornih kodova (216 autora)
Broj naredbi grananja i petlji	Zbroj svih pojavljivanja naredbi grananja i petlji(for, while, do, if, else if, else, switch)	$\ln(\text{zbroj}/\text{duljina})$	7
Broj ključnih riječi	Zbroj svih pojavljivanja ključnih riječi programskog jezika, konkretno C++.	$\ln(\text{zbroj_kljucne}/\text{duljina})$	1
Ternarni operatori	Broj pojavljivanja ternarnog operatora	$\ln(\text{broj_ternarnih}/\text{duljina})$	1
Komentari	Broj pojavljivanja komentara	$\ln(\text{broj_kom}/\text{duljina})$	1
Konstante	Broj pojavljivanja znakovnih i brojčanih konstanti	$\ln(\text{broj_konst}/\text{duljina})$	1
Makro naredbe	Broj pojavljivanja makro naredbi	$\ln(\text{broj_makro}/\text{duljina})$	1
Funkcije	Broj funkcija	$\ln(\text{broj_fun}/\text{duljina})$	1
Tokeni	Token je ekvivalentan unigramu, zbroj svih tokena	$\ln(\text{broj_token}/\text{duljina})$	1
Mjere duljine linija	Standardna devijacija i prosjek duljine linija	$\text{stddev}(\text{linije})$, $\text{avg}(\text{linije})$	2
Mjere parametara funkcija	Standardna devijacija i prosjek broja parametara funkcija	$\text{stddev}(\text{param})$, $\text{avg}(\text{param})$	2
Operatori	Zbroj pojavljivanja svih operatora programskog jezika	$\ln(\text{zbroj_op}/\text{duljina})$	1

1.1.2 Strukturne značajke

Strukturne značajke opisuju kakvu strukturu autor koristi dok piše izvorni kod, npr. koristi li tabove ili razmake na početku linije, piše li novu liniju prije nego otvori kontrolni blok, itd. *Tablica 1.2* detaljno opisuje svaku od korištenih značajki.

Tablica 1.2: Definicija strukturnih značajki

<i>Ime značajke</i>	<i>Definicija</i>	<i>Izraz</i>	<i>Veličina</i>
Tabovi	Broj svih tabova	$\ln(\text{broj_tabova}/\text{duljina})$	1
Razmaci	Broj svih razmaka	$\ln(\text{broj_razmaka}/\text{duljina})$	1
Omjer razmaka	Omjer razmaka(tabovi se broje) i ostalih znakova	$\ln(\text{zbroj_kljucne}/\text{duljina})$	1
Nova linija prije vitičastih zagrada	Koristi li autor novi red kada otvara vitičastu zagradu	boolean	1
Tabovi ili razmaci na početku linije	Koristi li autor tabove ili razmake na početku linije	boolean	1

Tablica 1.3: Definicija sintaksnih značajki

<i>Ime značajke</i>	<i>Definicija</i>	<i>Izraz</i>	<i>Veličina</i>
Bigrami čvorova	Bigrami čvorova su dva čvora koja su povezana u ASS-u	bigrami_čvorovaTF	dinamičko, oko 150000 za 2160 izvornih kodova
Tip čvora ASS-a	Frekvencija pojavljivanja tipa čvora ASS-a	tip_čvoraTF	58
Vrijednost lista ASS-a	Frekvencija vrijednosti listova ASS-a	list_vTF	dinamičko, oko 10000 za 2160 izvornih kodova

1.1.3 Sintaksne značajke

Sintaksne značajke se dobiju kako je već spomenuto iz apstraktnog sintaksnog stabla izvornog koda. One su što se vremena tiče najskuplje jer kreacija apstraktnih sintaksnih stabala nije brza, no trebale bi dati odlične značajke koje bi uvelike pomogle u deanonimizaciji. Apstraktna sintakсна stabla(ASS) su kreirana koristeći alat *joern* [2]. Ovaj alat nudi posebnu skriptu *joern-parse* koja parsira i vraća čvorove i bridove apstraktnog sintaksnog stabla. Postoji 58 različitih tipova čvorova(detilji u dodatku A) koje definira *joern.Tablica 1.3* detaljno opisuje svaku od korištenih značajki.

1.2 Selekcija značajki

Ovako kreirane značajke rezultiraju u ogromnim, rijetkim vektorima, čija veličina nekada doseže i stotine tisuća brojeva. Razlog tomu leži u definiciji značajki poput frekvencije tokena, frekvencije bigrama, itd. Rijetkost očitujemo u velikom broju nula unutar vektora. Rijetkost, također, može uzrokovati loš izabir idućeg čvora u klasifikatoru slučajne šume te s time i lošije rezultate. S velikim vektorima također dolazi i do puno sporijeg učenja klasifikatora jer će sva stabla odluka unutar slučajne šume imati više čvorova. Zbog svih navedenih razloga prije samog učenja klasifikatora napravljena je selekcija značajki koja odabire manji broj značajki koje sadrže dovoljno informacija da bi se klasifikator bolje i brže naučio. Tehnika selekcije značajki je mnogo pa ih ovdje neću sve detaljno opisivati nego ću se bazirati na tehnikama koje su korištene za ovaj diplomski rad, a to su selekcija značajki po

sadržaju informacije te selekcija značajki po iznosu varijance. Više o rezultatima sa i bez selekcije značajki u poglavlju 1.5.

1.2.1 Selekcija značajki po sadržaju informacije

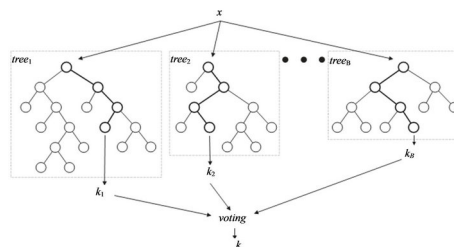
Svaka pojedinačna značajka vektora značajki nosi sa sobom neku količinu ili sadržaj informacije(eng. *information gain*), te nam ona igovori koliko je ta značajka bitna. Selektiramo samo one značajke koje sa sobom nose najveći sadržaj informacije. Definicija sadržaja informacije je detaljnije opisana u 1.3.2.

1.2.2 Selekcija značajki po iznosu varijance

Značajke su brojevi pa nad njima možemo računati razne statistike pa tako i varijancu. Ova selekcija značajki odbacuje sve značajke koje ne pređu unaprijed određenu granicu(eng. *threshold*) iznosa varijance.

1.3 Slučajna šuma

Slučajna šuma [3] je klasifikator koji se sastoji od kolekcije nezavisnih stabala odlučivanja. Svako od stabala predstavlja jedan glas u većinskom donošenju odluke. Odluka se donosi zbrajanjem glasova te se odabire odluka s najvećim brojem glasova [4]. Ovo detaljnije možemo vidjeti na slici 1.1. Slučajna šuma jer je u svojoj osnovi samo skup stabala vrlo dobro podnosi veliku dimenzionalnost podataka (što za ovaj problem očekujemo) i ne očekuje linearnu odvojivost vektora značajki te je iz tih razloga odabrana kao korišteni algoritam.



Slika 1.1: Arhitektura slučajne šume [5]

Svako od N stabala odluke je izgrađeno nasumičnim uzorkovanjem s ponavlja-

njima skupa za treniranje tako da se uzorkuje podskup duljine N . Stabla se grade do maksimalne moguće dubine iako postoje instance algoritma u kojem se stabla podrezuju. U izgradnji stabla ponovno se slučajno odabire podskup značajki kojih ima M . Veličina tog podskupa je hiperparametar algoritma, u literaturi [6] se za klasifikacijski problem preporuča veličina od \sqrt{M} . Od tog podskupa treba odabrati najbolju značajku koja će biti iskorištena za idući čvor stabla. Odabir najbolje značajke uobičajeno se radi metodama Gini nečistoće ili uzajamnog sadržaja informacije.

1.3.1 Gini nečistoća

Gini nečistoća je mjera koliko često bi nasumično odabrana značajka iz nekog skupa bila krivo klasificirana ako bi ju se nasumično klasificiralo s obzirom na to kakva je razdioba značajki po razredima u podskupu svih značajki. Drugim riječima gini nečistoća je kriterij koji teži minimizaciji vjerojatnosti krive klasifikacije [7]. Računamo ju na sljedeći način [8]:

$$I_g(t) = 1 - \sum_{i=1}^c p(i|t)^2 \quad (1.1)$$

gdje je $p(i|t)$ broj značajki koje pripadaju klasi i za čvor t .

1.3.2 Uzajamni sadržaj informacije

Uzajamni sadržaj informacije je koncept baziran na entropiji. Entropija je definirana kao količina informacije koju nosi neka poruka te ju računamo:

$$H(t) = - \sum_i p(x_i) * \log_2 p(x_i) \quad (1.2)$$

gdje su $p(x_i)$ vjerojatnosti svake od klasa.

Uzajamni sadržaj informacije definiran je kao:

$$I(X; Y_i) = H(X) - H(X|Y_i) \quad (1.3)$$

gdje je X klasa(autor), a Y_i i-ta značajka iz skupa. Intuitivno ga možemo zamisliti kao količinu informacije koju daje značajka i za klasu kojoj pripada.

1.4 Prikupljanje podataka

Izvorni kodovi korišteni u eksperimentima djelo su učenika srednjih i osnovnih škola koji su se natjecali na HONI-u¹ u godini 2016-2017. Skupljena su dva skupa podataka, jedan od 216 autora gdje svaki od njih ima 10 izvornih kodova i drugi od 29 autora s također 10 izvornih kodova, ali ovaj put su to izvorni kodovi kao rješenja istih 10 zadataka dok su u prvom primjeru oni birani nasumično.

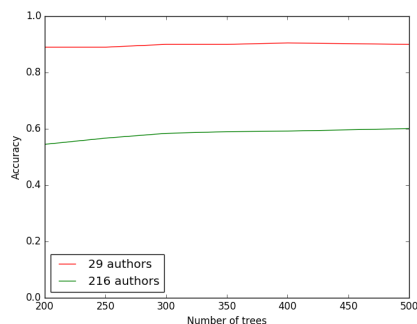
1.5 Rezultati i rasprava

Napravljeno je nekoliko različitih eksperimenata s različitim hiperparametrima i selekcijom značajki kako bi dobili što bolju sliku kako algoritam radi. Svaki od eksperimenata je odrađen k-preklop među validacijom (eng. *k-fold cross validation*), gdje je k u našem slučaju 10 jer imamo po 10 rješenja za svakog autora, kako bi dobili što objektivnije i točnije rezultate. U svakoj iteraciji 9 izvornih kodova autora birano je u skup za treniranje, a preostali izvorni kod u skup za testiranje. Također su prikazane po dvije slike za svaki eksperiment, jedna sa skupom od leksičkih i strukturnih značajki te jedna sa skupom leksičkih, strukturnih i sintaksnih značajki. To je napravljeno kako bi što bolje shvatili koliko su koje značajke bitne za prepoznavanje autora.

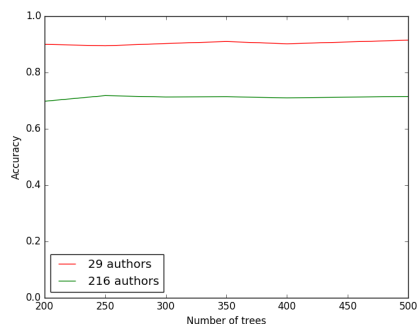
1.5.1 Rezultati bez selekcije značajki

U ovom poglavlju predstavljam rezultate bez selekcije značajki. Krenuli smo s pretpostavkom da će nam sintaksne značajke donijeti veliko poboljšanje klasifikatora što spominju i [1], no kao što možemo vidjeti na slikama 1.2 i 1.3 nad naša dva skupa podataka gotov jednaku točnost imamo u slučaju sa 29 autora, dok su na skupu sa 216 autora ukupnu točnost dosta smanjili.

¹<http://www.hsin.hr/honi/>



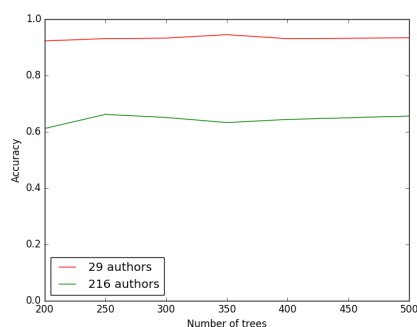
Slika 1.2: Točnost klasifikatora uz leksičke, strukturne i sintaksne značajke



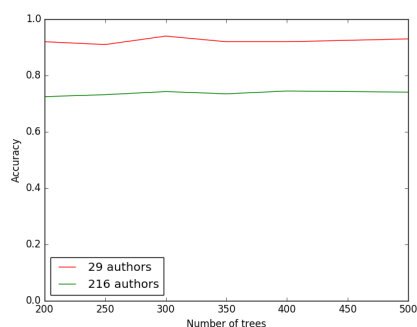
Slika 1.3: Točnost klasifikatora uz leksičke i strukturne značajke

1.5.2 Rezultati uz selekciju značajki sadržajem informacije

U ovom poglavlju predstavljam rezultate kada je korištena selekcija značajki uzajamni sadržajem informacije koja bi zbog prirode značajki koje su dosta rijetke(eng. *sparse*) trebala povećati točnost klasifikatora s obzirom na slučaj bez selekcije značajki što se doista i pokazalo istinito te to možemo vidjeti na slikama 1.4 i 1.5. Primjećujemo da su nam sintaksne značajke i u ovom slučaju donijele pogoršanje točnosti. Također pošto selekcija značajki smanjuje dimenzionalost vektora značajki sa stotina tisuća na desetke tisuća donosi nam veliko ubrzanje učenja klasifikatora.



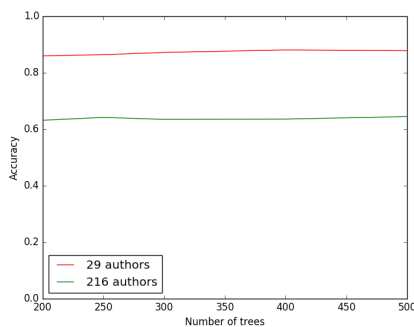
Slika 1.4: Točnost klasifikatora uz leksičke, strukturne i sintaksne značajke



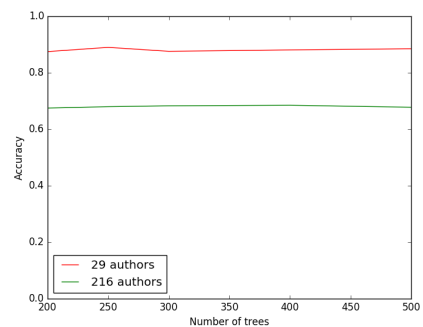
Slika 1.5: Točnost klasifikatora uz leksičke i strukturne značajke

1.5.3 Rezultati uz selekciju značajki varijancom

U ovom poglavlju predstavljam rezultate korištenjem selekcije značajki varijancom, u eksperimentu je korištena granica varijance od 0.03 jer je pokazala najbolje rezultate. Možemo primjetiti da su rezultati bolji nego kada smo koristili sve značajke, no lošiji su nego sa selekcijom uzajamnim sadržajem informacije što vidimo na slikama 1.6 i 1.7.



Slika 1.6: Točnost klasifikatora uz leksičke, strukturne i sintaksne značajke



Slika 1.7: Točnost klasifikatora uz leksičke i strukturne značajke

1.5.4 Rasprava

Prva stvar koju možemo primjetiti da što manji broj autora predamo klasifikatoru na treniranje to su nam rezultati bolji, što ima i smisla. Sintaksne značajke nam uglavnom donose lošije rezultate nego što smo pretpostavljali te vidjeli u [1]. Iz tog razloga te što nam donose stotine tisuća značajki i znatno usporavaju treniranje odlučio sam ih izbaciti za potrebe web aplikacije *Turtle*. Kao zaključak naveo bih da su rezultati dovoljno dobri, ali ipak nisam došao ni blizu rekreiranju rezultata iz [1].

1.6 Implementacijski detalji

Za implementaciju korišten je programski jezik *Python 2.7* te njegova biblioteka *scikit-learn* ² koja nudi pristup mnogim algoritmima strojnog učenja na vrlo jednostavan način. Kao implementaciju slučajne šume korišten je *RandomForestClassifier* ³. Selekciju značajki uzajamni sadržajem informacije implementirao sam pomoću *ExtraTreeClassifier* ⁴ koji omogućava biranje idućih čvorova pomoću sadržaja informacije te nam nakon učenja omogućava pristup najbitnijim značajkama. Selekciju značajki iznosom varijance implementirao sam korištenjem *VarianceThreshold* ⁵.

²<http://scikit-learn.org/stable/>

³<http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

⁴<http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.ExtraTreesClassifier.html>

⁵http://scikit-learn.org/stable/modules/generated/sklearn.feature_selection.VarianceThreshold.html

Određivanje sličnosti izvornih kodova

Ne postoji sustav specijaliziran za detekciju plagijata koji može sa sto postotnom sigurnošću utvrditi da je nešto plagijat. Kako bi odredili plagijat potreban je ljudski faktor. Odmah možemo uočiti da to nije baš uvijek efikasno, kada bi morali pronaći plagijate među tisućama dokumenata čovjeku bi trebalo puno vremena. Upravo iz tog razloga razvijamo sustav koji bi odredio sličnost među parovima dokumenata te izbacio parove za koje smatra da nikako ne mogu biti plagijat te uvelike ubrzao i olakšao posao ljudima. Dokumenti mogu biti teksutalne datoteke, izvorni kodovi, pjesme, itd.

Određivanje sličnosti izvornih kodova u svrhu detekcije plagijata je relativno neistraženo područje. Postoje dva vrlo slična, ali sada već stara sustava (nastali su prije više od 10 godina op.a) [9] [10] koji se baziraju na računanju otiska(eng. *documentfingerprint*) izvornog koda algoritmom *winnowing* koji je detaljnije opisan u [11]. *Turtle*, kao i sustav implementiran na završnom radu [12], koristi *winnowing* za račun otiska izvornog koda te preoblikovanje izvornog koda prije nego se otisak računa(detalnije opisano u radu [13]), no s dodatnim poboljšanjima. Neka od poboljšanja su veći broj preoblikovanja, što rezultira u boljim mjerama sličnosti i brže ukupno vrijeme. U nastavku poglavlja detaljnije su opisani koraci kojim sustav dolazi do mjera sličnosti, najvažnije pomoćne strukture podataka i algoritmi te usporedba s performansama starog sustava.

2.1 Izračun otiska izvornog koda

Otisak izvornog koda definiran je kao jedan podskup iz skupa izračunatih sažetaka (eng. *hash*) k-grama [11]. K-gram je uzastopni podniz duljine k . Odmah primjećujemo da k-grama ima $n - k + 1$ za znakovni niz duljine n te su ta dva broja vrlo blizu za manje vrijednosti k . Iako nigdje u literaturu nije navedeno koliko točno bi trebao biti k , najbolje rezultate su pokazale manje vrijednosti (otprilike prosjek duljine prosječnih duljina ključnih riječi programskog jezika na kojem se radi). Iz tog razloga potreban nam je efikasan algoritam kako bi izračunali sažetke svakog k-grama. Jedan od takvih algoritama je *Rabin-Karp* [14].

2.1.1 Rabin-Karp algoritam

Rabin-Karp je algoritam koji nam u linearnoj $O(n)$ složenosti računa sažetke svih k-grama ulaznog znakovnog niza. Algoritam radi tako da postupno gradi rješenje, iz i^{tog} sažetka se računa $i + 1$. Ako k-gram zapišemo kao $c_1...c_k$ tada sažetak računamo:

$$H(c_1...c_k) = c_1 * b^{k-1} + c_2 * b^{k-2} + ... + c_k \quad (2.1)$$

gdje su b baza, a c_i ascii vrijednost i^{tog} znaka. Za bazu se uzima prost broj kako bi se izbjegle kolizije.

Idući sažetak tada računamo tako da prethodno izračunatom sažetku oduzmemo prvi član, pomnožimo sve s b i dodamo novi znak:

$$H(c_2...c_k + 1) = (H(c_1...c_k) - c_1 * b^{k-1}) * b + c_{k+1} \quad (2.2)$$

Svaki uzastopni sažetak izračunat je uz samo tri operacije i tu možemo shvatiti linearnu složenost algoritma. Bitno je naglasiti da se iz ulaznog znakovnog niza prije računanja izbacuju nebitni znakovni poput razmaka.

2.1.2 Winnowing

Winnowing je algoritam za računanje otiska izvornog koda. Na ulaz prima niz sažetaka ranije izračunatih *Rabin-Karp* algoritmom te ispituje izračunati otisak. U nastavku je pseudokod algoritma:

Algorithm 1 Winnowing

Ulaz: K – niz sažetaka k -grama.

Izlaz: otisak F .

$\text{windows} := \text{split_to_windows}(w)$

for ($i := 0; i < n; ++i$) **do**

$w_min := \min(\text{windows}[i])$

if ($\text{more_than_one_same_val}$) **then**

if (w_min is rightmost) **then**

$F.\text{append}(w_min)$

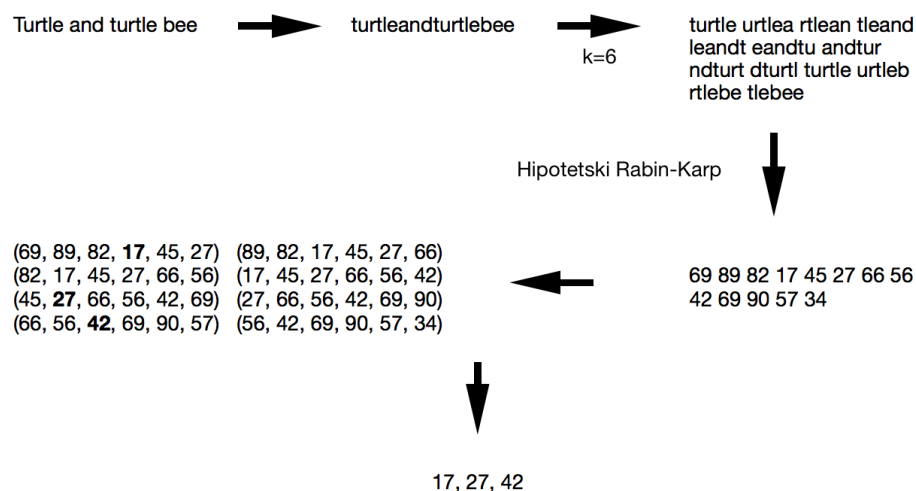
end if

end if

end for

return F

Primjer rada *Rabin-Karp* i *Winnowing* algoritama vidimo na slici 2.1.



Slika 2.1: Primjer rada algoritma *winnowing*

2.2 Preoblikovanje izvornog koda

Glavni razlog za korištenje preoblikovanja koda prije nego krenemo računati otisak je taj što algoritmi za računanje otiska pa tako i *winnowing* nisu otporni na mnoge pokušaje izmjene koda kako ih se ne bi detektiralo kao plagijate. Nabrojimo neke od njih:

- Dodavanje beskorisnih komentara
- Dodavanje beskorisnog koda
- Dodavanje praznina
- Zamjena imena varijabli
- Promjena jedne petlje u drugu
- Rastav naredbi na podbaredbe
- Izdvajanje koda u funkcije
- Promjena redoslijeda naredbi

Odmah možemo primjetiti da je *winnowing* otporan na dodavanje praznina jer se one brišu prije računanja sažetaka. Također izdvajanje koda u funkcije se detektira jer nema razlike gdje točno se isjecci koda nalaze za rezultat algoritma. Za ostale modifikacije trebamo napisati posebne funkcije koje će preoblikovati izvorni kod kako bi i one mogle biti detektirane. Naravno, modifikacije kao što je rastav naredbi na podnaredbe je jako teško detektirati bez korištenja jezičnog prevodica i neke vrste sintaksnih stabala te se tu limitiramo i njih *Turtle* ne zna detektirati. U nastavku poglavlja opisana su korištena preoblikovanja izvornog koda. Treba naglasiti da ovakvih preoblikovanja može biti jako puno te je sustav dizajniran tako da se svako preoblikovanje može vrlo lako dodati.

2.2.1 Zamjena imena varijabli

Ideja je da se sva imena varijabli zamjene univerzalnim imenom, u ovom slučaju odabrano je ime *var*. Također mijenjamo i sve konstante kako bi mogli detektirati i

promjenu redoslijeda naredbi. Prvo sve te varijable moramo pronaći unutar izvornog koda, a to radimo tako da izgradimo automat za željeni programski jezik koji će nam znati reći kada smo u stanju imena varijable. Kako bi automat znao razlikovati ključnu riječ od varijable potrebno je izgraditi strukturu podataka koja će mu reći promatra li ključnu riječ jezika. Za te potrebe izgrađeno je prefiksno stablo [15] od svih ključnih riječi koje nam tu funkcionalnost omogućava u jako brzom vremenu, zbog male duljine ključnih riječi možemo reći u $O(1)$. Implementirani automat također prepoznaje stanje komentar te izbacuje sve komentare iz izvornog koda pa time rješavamo i modifikaciju dodavanjem komentara. U dodatku B B vidimo izgled izvornog koda prije i poslije izvršavanja automata.

2.2.2 Zamjena while petlji u for petlje

Ovdje je opisan postupak za programski jezik C++, no postupak se može uz male preinake primjeniti i na sve poznatije programske jezike. Prvo primjetimo oblik *while* petlje:

```
1  init; while(test){ body; post; }
```

te ga zapišimo kao *for* petlju:

```
1  for(init; test; post){ body; }
```

Sada samo pronađemo ovakve isječke unutar koda i zamijenimo ih, npr. korištenjem regeks operacija.

2.3 Izračun sličnosti izvornih kodova

Nakon što smo izračunali otisak izvornog koda još samo trebamo izračunati sličnosti svih parova. Najprije se izračuna mapa indeksa (eng. *index mapping*) koja za svaki otisak sprema sve izvorne kodove u kojem se on nalazi. Nakon toga računaju se dvije sličnosti, sličnost prvog koda prema drugom kodu x_{12} te sličnost drugog koda prema prvom kodu x_{21} i to tako da se prebroje pojavljivanja otisaka koji su im zajednički te se taj broj normalizira ukupnim brojem otisaka prvog člana para. Dvije sličnosti imamo kako bi što bolje opisali slučajeve kada imamo dva izvorna

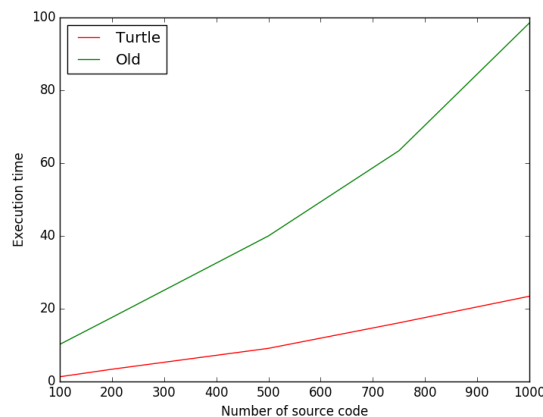
koda koja imaju veliku razliku u broju linija. S ovakvim računanjem sličnosti također rješavamo problem dodavanja beskorisnog koda prethodno opisanog. Za kraj nam samo ostaje definirati funkciju koja bi ove dvije sličnosti transformirala u jednu kako bi parove mogli sortirane prikazati u web aplikaciji:

$$f(x_{12}, x_{21}) = (x_{12} + x_{21}) * \min(x_{12}, x_{21}) \quad (2.3)$$

Odabrana je ova funkcija jer se u praksi pokazala kao odličan izbor. Ova funkcija nam također filtrira parove kako ne bi korisniku prikazivali parove koji sigurno nisu plagijati.

2.4 Usporedba sa starim sustavom

Kao što je spomenuto ranije ovaj sustav je novo implementirana verzija sustava implementiranog za završni rad te su u ovom poglavlju uspoređeni ti sustavi. Najveća prednost novog sustava je u njegovoj brzini što vidimo na slici 2.2.



Slika 2.2: Usporedba vremena izvršavanja

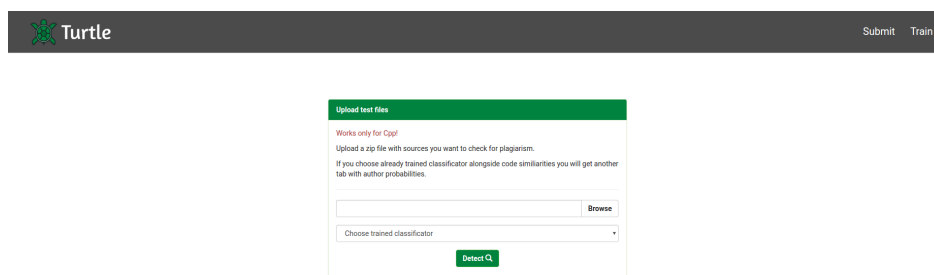
Mjereno je vrijeme od primitka izvornih kodova pa sve do ispisa parova sličnosti. Možemo primjetiti da oba sustava imaju linearnu složenost, no zbog puno bolje i pametnije implementacije *Turtle* je daleko bolji. *Turtle* također zbog više preoblikovanja izvornog koda nudi bolju detekciju nekih modifikacija (npr. zamjena *for* u *while* petlje). Ostatak modifikacija te izračunatih sličnosti nisu posebno opisivani jer su vrlo slični već opisanim u [12].

Integracija dva sustava


Svaka od komponenti sustava implementirana je u zasebnom projektu, sličnosti izvornih kodova su implementirane uz *Turtle* dok je određivanje autorstva implementirano u projektu nazvanom *Bee*. *Bee* je dizajniran kao skripta komandne linije kojoj se predaju skup za treniranje i skup za testiranje te koja vraća točnost naučenog klasifikatora. Pošto je *Turtle* web aplikacija bilo je potrebno nekako spojiti ova dva sustava. To je napravljeno tako da je u *Bee* napisan mali web server koji nudi nekoliko API endpoint-a za treniranje i predikciju. Te API endpoint-e koristi klijentski kod *Turtle* te su dvije aplikacije tako integrirane. Razloga za ovakvu integraciju je nekoliko, ali najvažniji je fleksibilnost, pošto sada funkcije *Bee* možemo koristiti kroz više različitih aplikacija.

3.1 Korisničko sučelje

U ovom potpoglavlju opisano je korisničko sučelje sustava *Turtle* kroz nekoliko slika.



Slika 3.1: Početna stranica, ovdje uploadamo zip datoteku koju želimo provjeriti, također možemo izabrati unaprijed trenirani klasifikator kako bi mogli dobiti i najvjerojatnije autore svakog koda.

 Turtle
 Submit Train

Upload train files

Name will be automatically updated when you upload your zip.
 Name your file uniquely so you can use it later at testing :)


Train

Trained files

All trained files will be shown here.
 User management is still not implemented!

train
train_names_chaged
train_files

Slika 3.2: Stranica za treniranje klasifikatora, uploada se zip datoteka te nakon što je učenje završeno naučeni klasifikator se pojavi u listi desno.

 Turtle
 Submit Train

Similarities
Deanonimization

#	First to second		Second to first		Detailed View
1	e14_033_kolo1_nepar.cpp e12_002_kolo1_nepar.cpp	85.45%	e12_002_kolo1_nepar.cpp e14_033_kolo1_nepar.cpp	87.61%	+
2	e03_013_kolo1_nepar.cpp e12_002_kolo1_nepar.cpp	100.00%	e12_002_kolo1_nepar.cpp e03_013_kolo1_nepar.cpp	80.53%	+
3	e03_013_kolo1_nepar.cpp e14_033_kolo1_nepar.cpp	100.00%	e14_033_kolo1_nepar.cpp e03_013_kolo1_nepar.cpp	74.55%	+
4	e12_002_kolo1_nepar.cpp e12_009_kolo1_nepar.cpp	53.98%	e12_009_kolo1_nepar.cpp e12_002_kolo1_nepar.cpp	46.83%	+
5	e03_013_kolo1_nepar.cpp e12_009_kolo1_nepar.cpp	67.03%	e12_009_kolo1_nepar.cpp e03_013_kolo1_nepar.cpp	41.27%	+

Slika 3.3: Nakon što sustav pronade sve sličnosti ispiše nam ih u tablicu sortirane od najveće do najmanje sličnosti.



Slika 3.4: Ako kliknemo na više detalja koji postoji za svaki par izvornih kodova dobijemo ovakav prozor gdje su obojani dijelovi koda koji su slični i što nam uvelike olakšava detekciju plagijata.

Turtle		Submit	Train
		Similarities	Deanonymization
e02_007_kolo5_funghi.cpp	e02_007	42%	
	e04_011	35%	
	e14_021	35%	
e02_021_kolo3_silveta.cpp	e02_021	42%	
	e505_001	35%	
	e29_003	35%	
e03_001_kolo4_gsenica.cpp	e03_001	42%	
	e03_014	35%	
	e14_008	35%	
e03_002_kolo6_voice.cpp	e03_002	42%	
	e14_044	35%	
	e12_010	35%	
e03_003_kolo6_voice.cpp	e03_003	42%	
	e14_058	35%	
	e101_008	35%	

Slika 3.5: Ako smo uz detekciju plagijata odredili da želimo i provjeriti tko bi mogli biti autori naših kodova na tabu *Deanonymization* dobijemo za svaki izvorni kod listu od 3 najvjerojatnija autora.

3.2 Implementacijski detalji

Turtle je implementiran u *Python2.7* koristeći njegov framework *Flask* ¹. Klijentski dio aplikacije implementiran je u *ReactJS* ² frameworku kako bi korisnici imali što bolji doživljan korištenja aplikacije te zbog asinkronog dohvata podataka.

¹<http://flask.pocoo.org/>

²<https://facebook.github.io/react/>

Zaključak

Ovaj sustav bio bi od velike pomoći velikom broju fakulteta, škola i brojnim organizatorima programerskih natjecanja jer u relativno kratkom roku može obraditi veliki broj podataka što je i dokazano testovima i uvelike bi smanjio pokušaje pisanja plagijata. Iako rezultati iz literature za deanonimizaciju autora nisu postignuti, opet su bili dovoljno dobri kako bi se mogli koristiti unutar sustava.

Duboko vjerujem da bi se ovaj sustav trebao nastaviti razvijati te kao najbitniju stavku proširenja vidim dodavanje novih programskih jezika poput *Pythona* i *Jave* jer se ti jezici se danas najviše koriste za pisanje laboratorijskih vježbi na fakultetima. Dalje kako bi bili otporniji na više modifikacija izvornog koda treba napisat dodatne postupke preoblikovanja koda. Trebalo bi smisliti i nove leksičke, strukturne i sintaksne značajke kako bi deanonimizacija autora radila bolje.

Kao završni cilj *Turtle* vidim kao sustav koji će fakultetima omogućiti lagano upravljanje laboratorijskim vježbama za niz predmeta koji bi uključivao sve od predaje rješenja do evaluacije istog te detekcije plagijata.

Literatura

- [1] A. Caliskan-Islam i dr. *De-anonymizing Programmers via Code Stylometry*. 2014. URL: https://www.princeton.edu/~aylinc/papers/caliskan-islam_deanonymizing.pdf.
- [2] F. Yamaguchi. *Joern documentation*. 2017. URL: <http://www.mlsec.org/joern/>.
- [3] L. Breiman. *Random Forests*. Machine Learning 45. 2001.
- [4] Nikola Bogunović. *Algoritmi strojnog učenja - 1, Strojno učenje*. predavanja, http://www.zemris.fer.hr/predmeti/kdisc/Algoritmi_1.ppt. Sveučilište u Zagrebu, Fakultet elektrotehnike i računarstva, Hrvatska.
- [5] C. Nguyen, Y. Wang i H. M. Nguyen. *Random forest classifier combined with feature selection for breast cancer diagnosis and prognostic*.
- [6] T. Hastie, R. Tibshirani i J. Friedman. *The elements of statistical learning*. 2009.
- [7] *DecisionTree*. predavanja, <http://www.cse.msu.edu/~cse802/DecisionTrees.pdf>. Michigan State University, USA.
- [8] Sebastian Raschka. <https://sebastianraschka.com/faq/docs/decision-tree-binary.html>.
- [9] Aiken A. URL: <https://theory.stanford.edu/~aiken/moss/>.
- [10] Karlsruhe Institute of Technology. URL: <https://jplag.ipd.kit.edu/>.
- [11] S. Schleimer, D. S. Wilkerson i A. Aiken. *Winnowing: Local Algorithms for Document Fingerprinting*. URL: <http://theory.stanford.edu/~aiken/publications/papers/sigmod03.pdf>.
- [12] D. Rakipović. *Sustav za detekciju plagijata*. 2015.

- [13] D. Šulc. *Algoritmi i metode mjerenja sličnosti izvornog koda*. 2015.
- [14] R. M. Karp i M. O. Rabin. *Pattern-matching algorithms*. 1987.
- [15] Paul E. Black. *trie*. online, Dictionary of Algorithms and Data Structures. 2011. URL: <https://xlinux.nist.gov/dads/HTML/trie.html>.

Sustav za detekciju plagijata te određivanje autorstva izvornih kodova

Sažetak

Plagijati postaju sve veći problem u današnjem svijetu, zbog ogromnog rasta broja informacija na internetu vrlo se jednostavno može ukrasti tuđi rad te predstaviti ga kao svoj, tj. plagirati. U sklopu ovog diplomskog rada naglasak je bio na pronalasku plagijata među izvornim kodovima te smo tom problemu pristupili iz dva kuta, prvi je pronalazak sličnosti među parovima izvornih kodova korištenjem otisaka izvornih kodova dok je drugi određivanje autora izvornog koda metodama strojnog učenja. Na kraju su ova dva pristupa spojena u web aplikaciju *Turtle* koja vrlo uspješno i u kratkom roku pronalazi plagijate među velikim brojem izvornih kodova (npr. laboratorijske vježbe ili programerska natjecanja).

Ključne riječi: Strojno učenje, plagijati, deanonimizacija, detekcija plagijata, izvorni kod

System for plagiarism detection and deanonymization of source code

Abstract

Plagiarism is becoming bigger and bigger problem in today society, the reason behind it is that information available online is growing exponentially and it's easy to steal from other authors and present it like your work. This thesis tackles the problem of detecting source code plagiarism from two different approaches, one is computing similarities between source code using source code fingerprints and other is deanonymizing authors of source code using machine learning. These approaches are then connected to a web application *Turtle*. This application can detect plagiarism between large number of source code (e.g. laboratory assignments, programming competitions) rather well.

Keywords: Machine learning, plagiarism, deanonymization, plagiarism detection, source code

DODATAK A

Tipovi čvorova apstraktnog sintaksnog stabla

Tablica A.1: Tipovi čvorova apstraktnog sintaksnog stabla

Additive Expression	AndExpression	Argument	ArgumentList
ArrayIndexing	AssignmentExpr	BitAndExpression	BlockStarter
BreakStatement	Callee	CallExpression	CastExpression
CastTarget	CompoundStatement	Condition	ConditionalExpression
ContinueStatement	DoStatement	ElseStatement	EqualityExpression
ExclusiveOrExpression	Expression	ExpressionStatement	ForInit
ForStatement	FunctionDef	GotoStatement	Identifier
IdentifierDecl	IdentifierDeclStatemetn	IdentifierDeclType	IfStatement
IncDec	IncDecOp	InclusiveOrExpression	InitializerList
Label	MemberAccess	MultiplicativeExpression	OrExpression
Parameter	ParameterList	ParameterType	PrimaryExpression
PtrMemberAccess	RelationalExpression	ReturnStatement	ReturnType
ShiftExpression	Sizeof	SizeofExpr	SizeofOperand
Statement	SwitchStatemetn	UnaryExpression	UnaryOp
UnaryOperator	WhileStatement		

DODATAK B

Primjer izvornog koda prije i poslije automata

```
1 #include <cstdio>
2 #include <iostream>
3
4 using namespace std;
5
6
7
8 int main () {
9
10     int a=0;
11     int sol=0;
12
13     for (int i=0; i<4; i++) {
14         cin >>a;
15         sol+=min (a, 12-a);
16     }
17
18     cout <<sol;
19
20     return 0;
21 }
```

Listing B.1: Prije

```
1 using namespace std;
2
3
4
5 int var () {
6
7     int var=var;
8     int var=var;
9
10    for (int var=var; var<var; var++)
11        {
12            cin >>var;
13            var+=var (var, var-var);
14        }
15    cout <<var;
16
17    return var;
18 }
```

Listing B.2: Poslije