

Multi-label klasifikacija klauzula iz pravnih ugovora

Katarina Aleksić

Softversko inženjerstvo i informacione tehnologije
Fakultet tehničkih nauka, Univerzitet u Novom Sadu
Novi Sad, Srbija
katarina.aleksic97@gmail.com

Nikolina Batinić

Softversko inženjerstvo i informacione tehnologije
Fakultet tehničkih nauka, Univerzitet u Novom Sadu
Novi Sad, Srbija
nina.batinic@yahoo.com

Sažetak— Poslednjih godina, sve je veći fokus na primeni veštačke inteligencije, naročito NLP (eng. *Natural Language Processing*) tehnika, u rešavanju zadataka iz pravnog domena. Kroz ovaj rad su opisana rešenja za problem klasifikacije klauzula iz pravnih ugovora u jednu ili više kategorija. Zadatak je da se iz jednog pravnog dokumenta izdvoje sve klauzule koje pripadaju jednoj ili više kategorija, te da se naznači koje su te kategorije. To bi u velikoj meri olakšalo i smanjilo napore za pregled dugih i kompleksnih pravnih ugovora od strane ljudi. Pravnicima bi u tom slučaju bilo mnogo lakše da pronađu najvažnije stavke iz ugovora. Kroz rad je predloženo nekoliko algoritama dubokog učenja i tehnika za *embedding* u cilju klasifikacije teksta. Skup podataka koji se koristi u radu je ručno labeliran od pravnih eksperata iz *The Atticus Project*-a¹ i sadrži preko 12.000 labeliranih klauzula.

Gljučne reči—NLP; neuronske mreže; LSTM; BERT; klauzule; EDGAR;

I. UVOD

Primena veštačke inteligencije u pravnom domenu dobija sve veći značaj u poslednjih nekoliko godina. Iako su se mnogi pravници ranije opirali promenama, nakon velikog napretka u NLP-u (*Natural language processing*), sve je veća zainteresovanost, kako pravnika, tako i inženjera koji se bave veštačkom inteligencijom[1].

NLP može da omogući brojne benefite pravnoj informatici i sve više se primenjuje u[2]:

1. Pravnom istraživanju – pronalaženje informacija relevantnih za pravnu odluku,
2. Elektronskom otkrivanju - utvrđivanje relevantnosti dokumenata za zahtev za informacijama,
3. Pregled ugovora(eng. *contract-review*) - provera da li je ugovor potpun i izbegavanje rizika,
4. Automatizacija dokumenata - generiranje rutinskih pravnih dokumenata,
5. Pravnom savetu: korišćenje dijaloga za pitanja i odgovore za pružanje prilagođenih saveta.

Ovaj rad bavi se pregledom ugovora.

Automatizovani sistemi za pregled ugovora mogu se koristiti za pregled dokumenata koji su relativno standardizovani i predvidljivi u pogledu vrsta sadržaja koje sadrže. Proces uključuje razlaganje ugovora na njegove pojedinačne odredbe ili klauzule, a zatim procenjivanje svake od njih, bilo da bi se izdvojile ključne informacije ili da bi se uporedile sa nekim standardom[2].

Prilikom pregleda ugovora, pravници moraju često da temeljno prolaze kroz čitave ugovore koji mogu biti dugački od nekoliko strana, pa do čak nekoliko stotina strana, samo da bi ručno pronašli par značajnih klauzula koje se nalaze u tom ugovoru, šta je tim klauzulama određeno, da li neke klauzule nedostaju u ugovoru, i slično. Prema podacima Razvojne banke Saveta Evrope²(eng. *Council of Europe Development Bank*, CEB), mnoge advokatske firme potroše oko 50% vremena na pregledanje ugovora[3].

Pregled ugovora, osim što zahteva veliku količinu vremena, kao i neefikasno korišćenje veština pravnih stručnjaka, često je i veoma skup proces za klijenta. Umesto da se pravници koncentrišu na značajnije probleme, oni moraju da troše svoje dragoceno vreme na dosadne zadatke koji mogu biti automatizovani.

Iako već postoje sistemi bazirani na pravilima koji mogu donekle da olakšaju ovaj posao, sa napretkom NLP-a stvaraju se mnogo veće mogućnosti za poboljšanje ove oblasti, jer je u osnovi svih ovih dokumenata jezik, tekst[4].

Glavni problemi u automatizaciji pregledanja ugovora i izdavanja klauzula jeste to što su izrazi koji se koriste u ugovorima vezani za specifičan pravni domen. Zbog toga, neophodni su podaci koji su vezani za ovaj domen, ali danas, iako postoji neki napredak, i dalje ima jako malo podataka koji su javno dostupni za obučavanje modela, i to ne samo u oblasti pregleda ugovora, već generalno u pravnoj informatici. Glavni razlozi za to su pre svega što su često u pitanju osetljive informacije koje ni ne smeju da budu javno dostupne, ali i što anotiranje tih podataka zahteva ekspertsko znanje pravnih stručnjaka, pa samim tim iziskuje dosta vremena i novca.

¹ <https://www.atticusproject.ai.org/>

² <https://coebank.org/en/>

U 2. poglavlju opisani su radovi koji su povezani i značajni za ovaj rad. U sledećem poglavlju opisan je skup podataka korišten u radu, a potom metodologije korištene za preprocesiranje podataka, *embedding* i klasifikaciju. U 5. poglavlju prikazani su dobijeni rezultati, i u poslednjem poglavlju je zaključak rada.

II. POVEZANI RADOVI

Kada uzmemo u obzir da je NLP svoj nagli napredak doživeo u poslednjih par godina, kao i da još uvek ima veoma malo javno dostupnih podataka vezanih za pravnu informatiku, nije ni čudo što ne postoji preveliki broj radova koji se bave ovom temom, kao i to da su svi koji postoje skorijeg datuma. U ovom poglavlju će biti navedeni neki od njih.

A. LEGAL-BERT

BERT³(*Bidirectional Encoder Representations from Transformers*) predstavlja *state-of-the-art* u nekoliko NLP zadatka sa generičkim podacima. Međutim, pošto je za treniranje BERT-a korišćen generički korpus kao što su podaci sa Vikipedije⁴(*Wikipedia*), dečije knjige i slično, BERT često daje slabije rezultate kada su u pitanju specifični domeni. Zbog toga postoje mnoge verzije BERT-a koje su specijalizovane za određene domene, poput biomedicinskog, tako što je postojeći BERT pretrenira ili se trenira od nule, i takvi domenski specijalizovani modeli daju značajno bolje rezultate.

U ovom radu[4] predstavljen je prvi, i za sada jedini model BERT-a koji je specijalizovan za pravni domen. Ovaj model daje značajno bolje rezultate u odnosu na standardni BERT u pogledu različitih složenijih problema. Što se tiče *multi-label* klasifikacije, poboljšanje je čak 2.5%. Obučavan je na čak 12 GB različitog pravnog teksta na engleskom jeziku iz nekoliko oblasti (npr. zakonodavstvo, sudski sporovi, ugovori).

Korpus za treniranje LEGAL-BERT-a sastoji se iz:

- 116.062 pravnih dokumenta Evropske unije, javno dostupnih na EURLEKS⁵-u, repozitorijumu koji je pod upravom Ureda za publikacije Evropske Unije⁶(eng. *EU Publication Office*).
- 61.826 dokumenata britanskog zakonodavstva, javno dostupnih sa portala britanskog zakonodavstva⁷.
- 19.867 predmeta Evropskog suda pravde (eng. *European Court of Justice*, ECJ), takođe dostupnih na EURLEKS-u.
- 12.554 predmeta sa HUDOC⁸-a (*Human Rights Documentation*), repozitorijuma Evropskog suda za

ljudska prava⁹ (eng. *European Court of Human Rights*, ECHR).

- 164.141 predmet različitih sudova širom SAD-a, dostupnih na portalu *Project Law Access Project*¹⁰.
- 76.366 američkih ugovora sa EDGAR-a (*Electronic Data Gathering, Analysis, and Retrieval*), baze podataka Američke komisije za hartije od vrednosti¹¹ (eng. *US Securities and Exchange Commission*, SECOM)

U radu su predstavljene različite varijacije LEGAL-BERT-a, koje su i javno dostupne za korišćenje, i to su:

- CONTRACTS-BERT-BASE - pretreniran na američkim ugovorima, pogodan za prepoznavanje imenovanih entiteta (eng. *Named entity recognition*, NER)
- EURLEX-BERT-BASE – pretreniran na dokumentima sa EURLEKS repozitorijuma, pogodan za *multi-label* klasifikaciju,
- ECHR-BERT-BASE – pretreniran na slučajevima Evropskog suda pravde, pogodan za binarnu i *multi-label* klasifikaciju,
- LEGAL-BERT-BASE – pretreniran na celokupnom korpusu,
- LEGAL-BERT-SMALL – treniran od nule na celokupnom korpusu, čak 33% manji od BERT-BASE-a, sa približnim peromansama kao on, ali čak 4 puta brži.

B. Revizija ugovora pomoću transformer modela

The Atticus Project predstavlja neprofitnu organizaciju sa ciljem da ubrza efikasnost pregledanja ugovora primenom veštačke inteligencije. Zahvaljujući tom projektu, sada postoji labelirani skup podataka sa preko 13000 labela iz 510 ugovora preuzetih sa EDGAR-a. U ovom korpusu postoji 41 labela. Veliki značaj ovog skupa podataka ogleda se u tome da on trenutno predstavlja jedini veliki skup podataka koji postoji za ovaj problem pregledanja ugovora i specijalizovan je za NLP zahvaljujući tome što su klauzule u ugovorima ručno su anotirane od strane pravnih stručnjaka i više puta pažljivo proverene.

Glavni cilj rada[3] je primena transformer modela da se iz ugovora bojenjem izdvoje važne klauzule kako bi advokatima, ljudima, bio olakšan pregled ugovora.

Većina klauzula u ugovoru nije labelirana jer one nisu važne prilikom pregledanja ugovora, što dovodi do velike neravnoteže između relevantnih i nerelevantnih klauzula u

³ <https://github.com/google-research/bert>

⁴ <https://en.wikipedia.org/>

⁵ <https://eur-lex.europa.eu>

⁶ <https://op.europa.eu>

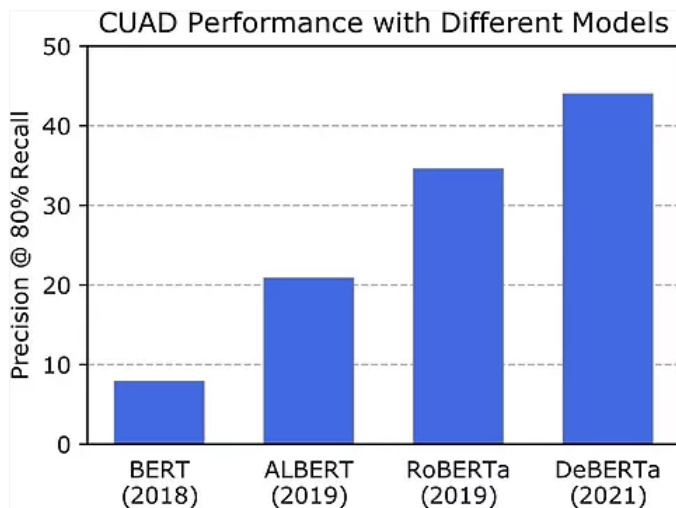
⁷ <https://www.legislation.gov.uk>

⁸ <https://hudoc.echr.coe.int/>

⁹ <https://www.echr.coe.int/>

¹⁰ <https://case.law>

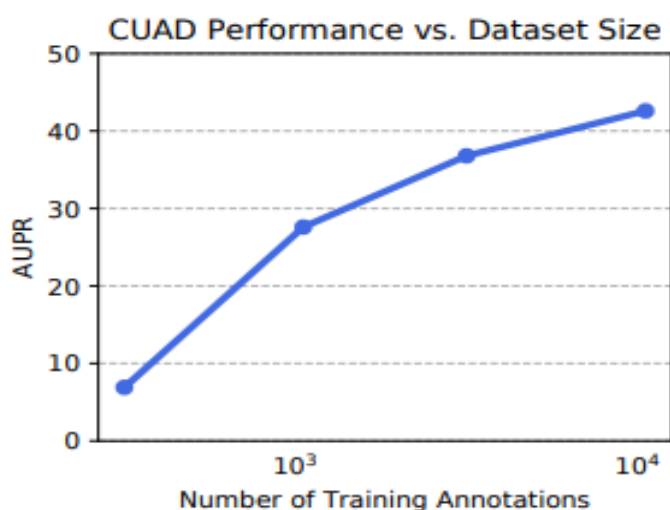
¹¹ <https://www.sec.gov/edgar.shtml>



Slika 1. Performanse različitih modela korišćenih u radu nad istim skupom podataka[6].

opozivu(eng. *Precision at recall*) mera performansi modela. Npr. u ugovorima postoji 100 relevantnih klauzula, a model izdvoji 500 klauzula, onda je preciznost modela $100/500=20\%$. Od tih 500, ako je 80 klauzula je relevantno, onda znači da je pronšao 80 od 100 relevantnih klauzula, pa je odziv 80%. Korisniku u ovom slučaju neće biti izdvojeno 20 važnih klauzula i moraće da pročita dodatnih 420 nerelevantnih klauzula. U ovom primeru mera je 20% *precision@80% recall*.

Korišćena mera je i veličina površine ispod krive koja meri preciznost pri opozivu(eng. *Area Under the Precision-Recall curve*, AUPR). S obzirom na to da svaka preikcija ima svoju verovatnoću pouzdanosti, pa se može postaviti prag(eng. *Threshold*) kolika je minimalna vrednost pouzdanosti neophodna za predikciju, i na osnovu toga dobijati željene opozive, a zatim za svaku vednost opoziva izračunati preciznost.



Slika 2. Performanse modela u zavisnosti od veličine skupa podataka za obučavanje[6].

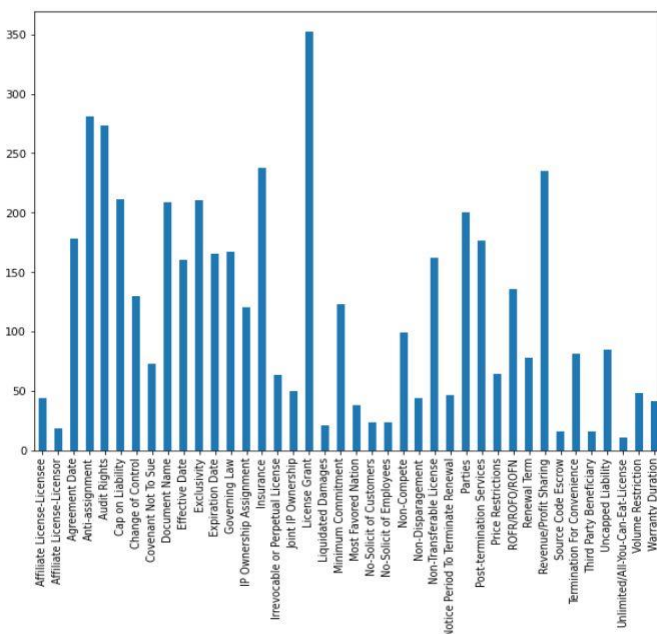
Modeli koji su korišćeni i njihove performanse su prikazani su na slici 1. Iz priloženog se može zaključiti da kako se poboljšavaju NLP modeli, tako se poboljšava i tačnost.

Međutim, mnogo veći uticaj na performanse ima veličina skupa podataka. Na slici 2 prikazan je nagli skok performansi kada broj podataka dođe do 10^4 , što nije ni čudo, s obzirom na to da su ovo kompleksni modeli i da se za obučavanje transformer modela koriste i milioni podataka.

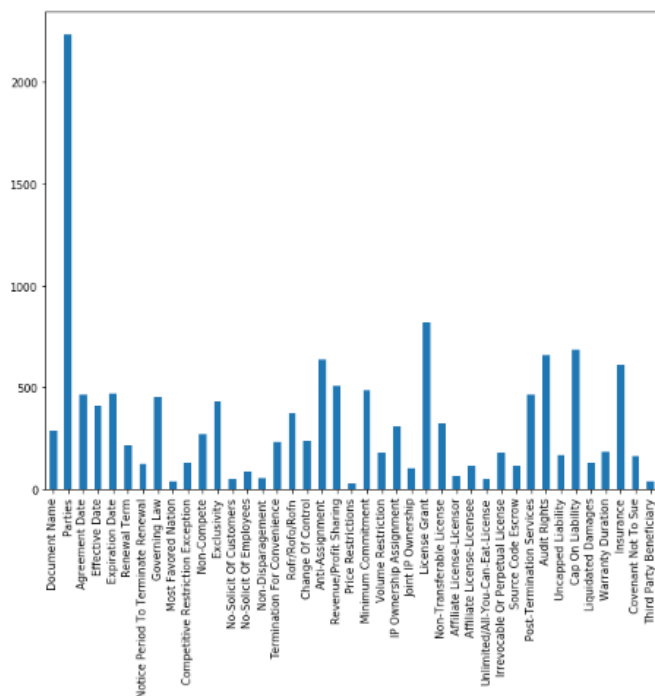
Projekat Atikus ne prestaje sa daljim razvojem. Skup podataka, modeli i rezultati konstantno se poboljšavaju.

III. SKUP PODATAKA

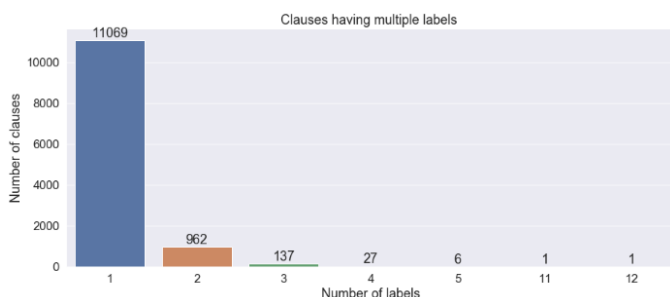
Na početku rada bila je dostupna verzija skupa podataka koja se sastojala se od 200 pravnih ugovora dostupnih u pdf i csv formatu. Za treniranje i testiranje modela korišteni su dokumenti koji sadrže labelirane klauzule u csv formatu. Svi ugovori su komercijalni i potiču iz sistema EDGAR koji koristi Komisija za hartije od vrednosti i berze SAD (SEC). Važna karakteristika ovih ugovora je da su kompleksniji i da sadrže veliki broj klauzula koje je teško naći u opštoj populaciji ugovora. Podaci su prikupljeni i labelirani od strane neprofitne organizacije *The Atticus Project* čiji je cilj da se iskoriste sve prednosti veštačke inteligencije u pravnom domenu, naročito u postupku pregleda pravnih ugovora. Svaki od komercijalnih ugovora spada u jedan od 25 tipova ugovora. Neki od tipova ugovora su: Sporazum o zajedničkom brendiranju (eng. *Co-Branding Agreement*), Ugovor o licenciranju (eng. *Licence Agreement*) itd. Iz tog razloga ugovori su prilično varirajućih dužina, od onih koji zauzimaju samo par strana, do kompleksnih ugovora sa preko 100 strana.



Slika 3. Broj primera koji pripadaju svakoj od 40 kategorija za prvu verziju skupa podataka.



Slika 4. Broj primera koji pripadaju svakoj od 41 kategorija za drugu verziju skupa podataka.



Slika 5. Koliko klauzula spada u koliko kategorija

Organizacija *The Atticus Project* se zasniva na volonterima, među kojima su i pravni stručnjaci. Upravo oni su kreirali i ručno labelirali skup podataka tako da se sastoji od preko 4000 klauzula kategorisanih u jednu ili više kategorija. Ukupno je 40 kategorija, koje iskusni pravници smatraju važnim u postupku pregledanja ugovora. Sve kategorije su međusobno nezavisne.

Neke od 40 kategorija su:

- *Governing Law* – klauzula u kojoj se navodi u kojoj državi se reguliše ugovor (primer: “*This Agreement is accepted by Company in the State of Nevada and shall be governed by and construed in accordance with the laws thereof, which laws shall prevail in the event of any conflict.*”);
- *License Grant* – klauzula u kojoj se navodi da li ugovor sadrži licencu koju je jedna od strana odobrila drugoj ugovornoj strani (primer: “*i-Escrow hereby grants to 2TheMart a worldwide, non-exclusive right*

to use, reproduce, distribute, publicly perform, publicly display and digitally perform the i-Escrow Content on or in conjunction with 2TheMart auctions.”);

- *Non-Disparagement* – klauzula kojoj se navodi da li postoji zahtev jedne ugovorne strane da druga strana ne sme da priča negativno o njoj. To je najčešće slučaj kod poslodavca i zaposlenog, gde zaposleni potpisuje da kompaniju u kojoj radi neće spominjati u negativnom kontekstu u bilo kojoj formi komunikacije (primer: “*The Company shall not tarnish or bring into disrepute the reputation of or goodwill associated with the Seller Licensed Trademarks or Arizona.*”);
- *Expiration Date* – klauzula u kojoj se navodi kada ističe prvobitni rok ugovora (primer: “*The term of this Agreement shall continue for one (1) year following the Launch Date, unless earlier terminated as provided herein.*”).

Dobijeni skup podataka je prilično nebalansiran, što se može videti na slici 3, zbog čega je isproban i *oversampling* podataka. Na primer, labela „*Liquidated Damages*“ sadrži samo 21 primer (toliko klauzula pripada toj kategoriji), dok labela „*License Grant*“ ima 352 primera. Takođe, podaci su nebalansirani i u smislu broja kategorija kojima labele pripadaju: 3751 klauzula pripada samo jednoj kategoriji, 422 labele pripadaju dve kategorije, 47 klauzula spadaju u 3 kategorije i samo 2 klauzule spadaju u 4 kategorije.

Ova verzija skupa podataka kasnije je dopunjena. U drugoj verziji postoji 41 labela, 500 ugovora i preko 13000 klauzula, ali je i dalje ostao problem nebalansiranih podataka, što se može videti na slici 4. Kategorija koja nije postojala u prvoj verziji je *Competitive Restriction Exception*. U ovom skupu postoji klauzula sa čak 12 kategorija. Na slici 5 može se videti koliko klauzula spada u koliko kategorija.

IV. METODOLOGIJA

A. Pretprocesiranje podataka

Priprema podataka (eng. *Data wrangling*) predstavlja proces prikupljanja, transformisanja i čišćenja podataka. Osim toga, mnogi drugi aspekti poput provere kvaliteta podataka, spajanja sa različitih izvora, poboljšanja sa drugim podacima i slično, predstavljaju bitan deo pripreme podataka[3]. Ovaj proces je krucijalan u postupku klasifikacije podataka, jer je to jedini način da sirovi podaci postanu upotrebljivi.

Za prvu verziju skupa podataka, prvi korak u pretprocesiranju podataka za *multi-label* klasifikaciju klauzula u pravnim ugovorima je prikupljanje podataka i spajanje labeliranih klauzula iz svakog csv dokumenta u jedan *dataframe*. Nakon toga smo izbacili nepotrebne kolone (npr. *Label 1-Answer*, *Label 2-Answer*). S obzirom da vrednosti kolona predstavljaju nazive kategorija, potrebno ih je transformisati u nove kolone, gde će za svaku klauzulu vrednost tog reda i kolone iznositi 1 ako klauzula pripada toj kategoriji, a 0 ako ne pripada. Nakon ovog postupka *dataframe* ima 41 kolonu, gde su vrednosti prve kolone same klauzule, dok

ostalnih 40 predstavljaju svaku od kategorija kojima te klauzule potencijalno pripadaju, tj. ciljne labele, čije su vrednosti 1 ili 0.

U drugoj verziji skupa podataka, podaci su formatirani na drugačiji način.

Za svaku od labele postoji poseban *xlsx* fajl koji predstavlja izveštaj o nekoj labeli, zbog čega je i naziv svakog od fajlova u formatu "*Label Report – labela*". Neki od tih fajlova sadrže više različitih labele, npr. fajl sa datumima sadrži sledeće labele: *Agreement date*, *Effective date*, *Expiration date*, *Renewal term*, *Notice period to terminate renewal*.

Svi fajlovi organizovani su tako da prvu kolonu čini naziv ugovora u pdf formatu, zatim naziv klauzule, dok fajl sa datumima i fajl sa strankama posle svake od klauzula imaju još i po kolonu sa odgovorom za svaku od labele.

Još jedan fajl koji je značajan je *master.csv*, čija je prva kolona naziv ugovora, potom kolona sa nazivom labele, pa kolona sa odgovorom, npr. kolonu *Exclusivity* prati kolona *Exclusivity-Answer*, i tako za sve labele. Neke labele nisu identično nazivane u *Label Report* fajlovima kao u *master.csv* fajlu. Npr. kolona za labelu *Rofr/Rofo/Rofn* u *master.csv* fajlu nazvana je *ROFR-ROFO-ROFN* u svom *Label Report* fajlu i njihovi nazivi su ručno menjani da budu identični na oba mesta jer je to od značaja za kasnije pretprocesiranje podataka.

Takođe, ukoliko u jednom ugovoru postoji više klauzula sa istom labelom, one se nalaze u jednom polju, s tim što se u fajlu *master.csv* svaka od klauzula nalazila se između jednostrukih ili dvostrukih navodnika, a više klauzule koje su u istom polju međusobno su odvojene zarezom. Dok se u *Label Report* fajlovima klauzule ne nalaze pod navodnicima, a posle svake klauzule postoji još i podatak na kojoj strani ugovora je smeštena ta klauzula, dok je više klauzula u istom polju odvojeno pomoću "*\n\n*".

Pretprocesiranje se prvo sastoji od učitavanja fajla *master.csv* u *dataframe*. Iz *dataframe*-a je kolona sa nazivom ugovora preimenovana je u kolonu *Clauses*, gde treba smestiti tekst klauzula koje će se dodavati. Sve kolone sa odgovorima se uklanjaju, tako da ostanu samo kolone koje sadrže nazive labele.

Sledeći problem bio je kako izdvojiti same klauzule. U fajlu *master.csv*, s obzirom na to da i sam tekst klauzula može sadržati i jednostruke i dvostruke navodnike, kao i zapete, pa ih nije ih bilo moguće izdvojiti na taj način. Zato su podaci iz *dataframe*-a izbrisani, ali je služio za kasnije dodavanje klauzula.

Same klauzule su izdvojene iz *Label Report* fajlova. Potrebno je učitati svaki od fajlova, potom obrisati kolonu sa nazivom ugovora, a ukoliko je fajl sa datumima ili strankama obrisati i kolone sa odgovorima koje samo oni sadrže. Zatim se prolazi kroz sve preostale kolone i za svako od polja u koloni uklanjaju se delovi koji govore na kojoj strani ugovora se nalazi klauzula pomoću regularnog izraza "*\\(Page \\d+\\)*", a zatim se *split*-uje po "*\n\n*". Tako razdvojene klauzule smeštene su u listu, pa je potrebno proći kroz svaku od tih pojedinačnih klauzula u listi, i ukoliko se ne nalazi već u *dataframe*-u napravljenom na osnovu *master.csv* fajla dodati novi red sa tekstom klauzule i nulama za ostale kolone, a zatim,

bilo da je već postojala ili ne, u *dataframe*-u se postavi 1 za datu klauzulu i datu labelu.

Ostale faze identične su za obe verzije skupa podataka.

U sledećoj fazi pretprocesiranja su najpre uklonjeni svi specijalni karakteri iz klauzula, zatim reči koje sadrže samo jedno slovo, višestruke prazne linije itd. Nakon toga je primenjen *stemming* algoritam. *Stemming* predstavlja postupak svodenja različitih gramatičkih oblika reči poput imenica, prideva, glagola itd. na njihov korenski oblik. Osnovna razlika između *stemming*-a i lematizacije je to što lematizacija vrši i morfološku analizu reči zbog čega je neophodno proslediti i informaciju o vrsti reči[6]. Završna faza pre *embedding*-a klauzula predstavlja izbacivanje zaustavnih reči (eng. *stop words*). *Stop words* predstavljaju reči koje se najčešće pojavljuju u svakom prirodnom jeziku. U svrhu analize teksta i formiranja NLP modela za za procesiranje teksta, ove reči najverovatnije neće dati nikakvu vrednost značenju dokumenta, zbog čega se najčešće izbacuju. Za klasifikaciju klauzula u PDF dokumentima, potrebno je najpre izvršiti tokenizaciju dokumenta (transformacija dokumenta u sekvencu rečenica).

Nakon pripreme i pretprocesiranja podataka prelazi se na fazu *word embedding*-a. *Word embedding* predstavlja transformaciju teksta u vektorski oblik koji kodira značenje reči tako da se očekuje da reči koje su slične po značenju imaju sličnu vektorsku reprezentaciju. Tehnike *embedding*-a korištene u ovom radu su sledeće: *GloVe* i *Bert*.

1) GloVe

GloVe predstavlja nenadgledan algoritam učenja za vektorsku reprezentaciju reči. Osnovna karakteristika ovog modela je da se oslanja na matricu brojeva međusobnog zajedničkog pojavljivanja reči (eng. *co-occurrence matrix*). Ta matrica se formira za okolinu (prozor oko) reči. Ideja na osnovu koje se formiraju *GloVe* vektori je - Rastavljanje *co-occurrence* matrice na delove pomoću metode koja se zove *Singular Value Decomposition* (SVD). S obzirom da se postupak primene SVD u praksi pokazao kao problematičan, autori *GloVe*-a su SVD sveli na optimizacioni problem, zbog čega je napravljen jako robusan i skalabilan model. Jednostavno rečeno, *GloVe* omogućava da nad korpusom reči transformišemo svaku reč u određenu poziciju u višedimenzionom prostoru. To znači da će reči sa sličnim značenjem biti relativno blizu.

Za vektorizaciju pravnih klauzula isprobali smo i pretrenirani *GloVe* vektor¹². On predstavlja korpus reči gde je svaka reč praćena sa 100 brojeva koji opisuju vektor pozicije reči. Najpre je nad klauzulama primenjen *tokenizer* koji kreira rečnik indeksa, gde je vrednost svake reči broj njenih pojavljivanja u svim klauzulama. Nakon toga je od rečnika indeksa formirana sekvencu brojeva, nakon čega svaka rečenica sadrži numeričku reprezentaciju reči. S obzirom da sekvence nisu jednakih dužine nad njima je izvršena *pad_sequence* metoda (vrednosti koje nedostaju se popunjavaju 0). Zatim se primenjuje pretrenirani *GloVe* vektor od koga se formira *embedding* matrica koja se kasnije će se kasnije koristiti kao vrednost početnih težina u *Embedding* sloju kod rekurentnih neuronskih mreža.

¹² <https://www.kaggle.com/danielwillgeorge/glove6b100dtx>

B. Arhitektura rešenja

Za klasifikaciju klauzula koristili smo modele dubokog učenja koji su dali različite rezultate. Korišteni modeli su: LSTM (eng. *Long Short-Term Memory*) i BERT.

1) BERT

Transformer model prvi put je predstavljen 2017. godine. Namenjen je razumevanju jezika (eng. *language understanding*) pomoću *attention* mehanizma koji se ubacuje između enkodera (eng. *encoder*) i dekodera (eng. *decoder*) i tako dekoderu obezbedi mnogo više informacija, a ne samo poslednje skriveno stanje enkodera.. Za razliku rekurentnih neuronski mreža (eng. *Reccurent Neural Network – RNN*) koje čitaju reči sa leva na desno (ili obrnuto), transformer modeli, već se sve reči šalju odjednom, a redosled reči u rečenici pamti zahvaljujući pozicionom embeddingu (eng. *positional embedding*).

Transformeri koriste *self-attention*, koji omogućava da odredi koliko svaka od reči u sekvenci utiče na trenutno posmatranu reč tako što joj se dodeljuje *attention* koeficijent. *Self-attention* zapravo predstavlja težinsku kombinaciju (engl. *weighted combination*) svih *word embedding*-a u datom kontekstu, uključujući i reči koje se pojavljuju pre i reči koje se pojavljuju posle u rečenici, a kao težine se uzimaju *attention* koeficijenti

BERT predstavlja bidirekcionni transformer model koji je nastao u Google-u kao primena transformer arhitekture 2018. godine. On je dostigao *state-of-the-art* performanse u mnogim NLP zadacima. Zasniva se na dve tehnike: maskiranje reči (eng. *Masked Language Model – MDM*), gde se maskira 15% nasumično izabranih tokena sa ulaza, i predikcija da li jedna rečenica sledi iza druge (eng. *Next Sentence Prediction – NSP*). BERT kodira ulaznu rečenicu, koja predstavlja ulaz, kroz tri tipa *embedding*-a :

- *Token embedding* (za svaku reč se dobija njen indeks u rečniku) – prvi token, koji se naziva CLS token, koristi se za klasifikaciju zajedno sa *softmax* slojem;
- *Sentence embedding* (za svaku reč se dodaje indeks rečenice kojoj pripada);
- *Positional embedding* (za svaku reč se dodaje njena relativna pozicija u rečenici).

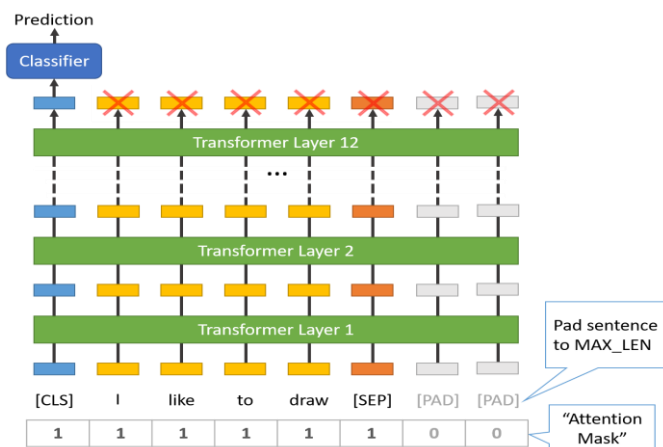
Sam BERT se sastoji od enkoder blokova (12 ili 24, u zavisnosti da li je u pitanju *Base* ili *Large* model)[9]. BERT ima ograničen vokabular na 30000 reči, tako da one reči koje ne pronalazi u rečniku razbija na n-grame.

U ovom radu korišten je pretrenirani LEGAL-BERT-BASE model o kome je bilo više reči u prethodnom poglavlju. On se sastoji iz 12 slojeva i ima 110 000 000 parametara.

LEGAL-BERT je *fine-tune*-ovan za *multi-label* klasifikaciju i to na sledeći način (slika 6):

- maksimalna dužina tokena postavljena je na 512, a ujedno je to i gornja granica koju BERT prihvata. Iako većina primera ima manje od 400 reči, 5 primera ima preko 512 tokena,

- stopa učenja (eng. *learning rate*) postavljen je na 0.00002, kako preporučuju autori BERT-a,
- veličina *batch*-a (*batch size*) postavljen je na 12, s obzirom na to da CUDA, platforma za paralelnu obradu na Nvidia¹³ grafičkom procesoru, nije imala dovoljno memorije za veći broj primera,
- broj epoha varirao je u zavisnosti od trening skupa. Bez *oversampling*-a, bilo je potrebno 10 epoha, dok su su *oversampling*-om bilo dovoljne 4 epohe,
- broj koraka po epohi dobija se kada se broj primera u trening skupu podeli sa veličinom *batch*-a i zaokruži nagore,
- ukupan broj koraka pri treningu jednak je proizvodu broju koraka po epohi i broja epoha,
- broj koraka za zagrevanje (eng. *warm up steps*) jednak je petini ukupnom broju koraka pri treningu,
- optimizator (eng. *optimizer*) – AdamW algoritam
- *learning rate scheduler* - *scheduler* koji u periodu zagrevanja linearno povećava stopu učenja od 0 do 0.00002, a nakon toga, do kraja treninga je linearno smanjuje na 0.
- klasifikator (eng. *classifier*) – za klasifikaciju se koristi linearni modul (eng. *linear module*) koji predstavlja *feed forward* sloj koji linearnom transformacijom konvertuje CLS token BERT-a (vektor bogat semantikom koji se sastoji od 768 float-ova koji predstavljaju *embedding* čitave rečenice) u 41 klasu, a potom se na taj dobijeni izlaz primeni *softmax* funkcija da bi se dobio broj u rasponu od 0 do 1.
- funkcija gubitka (eng. *loss function*) koja je korištena pri kompajliranju modela je *binary_crossentropy*. Ona se koristi za zadatke binarne klasifikacije (odgovor na pitanje ima dve opcije – 0 ili 1) i bitna karakteristika je da može istovremeno prediktovati podatke za više klasa, što je pogodno za *multi-label* klasifikaciju

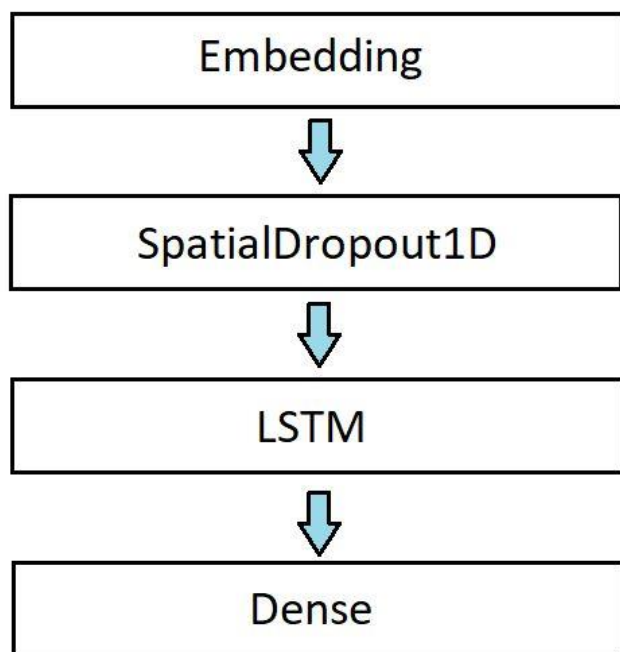


Slika 6. *Fine-tune*-ovan BERT za klasifikaciju.

¹³ <https://www.nvidia.com/>

2) LSTM

Rekurentne neuronske mreže predstavljaju jednu od najpopularnijih arhitektura koje se koriste u NLP-u. One su vrsta neuronskih mreža gde se izlaz neurona vraća na njegov ulaz, čime se pravi neka vrsta internog stanja mreže – memorije. LSTM predstavlja specijalnu arhitekturu RNN-a, uz pomoć koje je rešen problem kratkoročne memorije kod rekurentnih mreža. Uvode se tri nove komponente: *Input gate*, *Forget gate* i *Output gate*. Sve komponente predstavljaju obične logističke sigmoidalne funkcije, zbog čega interno stanje (memorija) nije vidljivo spolja, već samo kroz aktivacionu funkciju.



Slika 7. Arhitektura LSTM mreže korištene za klasifikaciju klauzula u pravnim ugovorima.

Kao što se može videti na slici 7, arhitektura LSTM mreže korištene u ovom radu se sastoji od 4 sloja. U pitanju je sekvencijalni model sa *Embedding*-om kao prvim slojem, koji omogućava konverziju svake reči u vektor fiksne veličine. Rezultujući vektor ima stvarne vrednosti, umesto 0 ili 1. Za težine mu je prosleđena prethodno formirana *embedding* matrica uz pomoć *GloVe* pre-treniranog vektora. Parametar *trainable* je *False*, iz razloga što koristimo pre-trenirani *word embedding*. Dimenzija ulaza (*input_dim*) predstavlja veličinu vokabulara, što je u ovom slučaju oko 4000, dok je dimenzija izlaza (*output_dim*) 100. Nad izlazom je potom primenjen *SpatialDropout1D* sloj, kako bi se sprečilo preprilagođavanje (eng. *overfitting*), čija vrednost parametra *rate* iznosi 0.2.

LSTM sloj je bidirekcionni sloj, što znači da model obrađuje sekvencu od početka do kraja, kao i unazad. Ovo je veoma bitno za razumevanje konteksta reči u rečenici. *Dense* sloj predstavlja potpuno povezani izlazni sloj sa 41 izlaznim neuronom, gde je vrednost svakog izlaza verovatnoća da

klauzula pripada labeli na toj poziciji. Sigmoid aktivaciona funkcija vraća vrednost 1 ili 0. Ako je vrednost svakog izlaznog neurona veća od 0.5, smatra se da klauzula pripada labeli koju predstavlja taj neuron.

Funkcija gubitka je *binary crossentropy*. *Optimizer* korišten pri kompajliranju modela je *Adam*. Model je treniran u 30 epoha, a *batch size* iznosi 64.

V. EVALUACIJA I REZULTATI

Skup podataka, koji se sastoji od 12203 primera, podeljen je na trening i test u razmeri 80:20. Za optimizaciju parametara upotrebljeno je 20% trening skupa za validaciju, koristeći *PredefinedSplit* validator (princip indeksiranja – podatak koji spada u trening set ima indeks -1, dok podatak koji spada u validacioni set ima indeks 0). Iako se unakrsna validacija često primenjuje kada je skup podataka nad kojim se model trenira i testira relativno mali, u ovom slučaju koristili smo poseban validacioni skup za optimizaciju modela zbog hardverskih ograničenja (treniranje modela je vremenski zahtevno).

S obzirom na veliku razliku u broju primera različitih klasa, isprobani su različiti načini balansiranja skupa podataka. Uklanjanje primera većinskih klasa (eng. *undersampling*) u ovom radu nije dobro rešenje s obzirom na to da manjinska klasa (eng. *minority class*) *Price Restrictions* ima svega 28 primera. Zbog toga su dodati primeri manjinskih klasa (eng. *oversampling*). *Oversampling* je izvršen nad skupom podataka za trening korištenjem *RandomOverSampler*-a koji na slučajan način bira primere manjinskih klasa koje će kopirati, i to na 2 načina:

1. Dodati su primeri za sve klase, osim *Parties*, tako da sve klase imaju isti broj primera kao i ona, s obzirom na to da *Parties* predstavlja klasu sa najvećim brojem primera (eng. *majority class*) i ima ukupno 2234 primera.
2. Dodati su primeri za sve klase, osim *Parties*, tako da sve ostale klase imaju duplo manje primera u odnosu na nju.

Razlog za ovo je sledeća klasa sa najvećim brojem primera nakon klase *Parties* jeste *License Grant* i ona ima 822 što je više nego duplo manje. Takođe, BERT istreniran na nebalansiranim podacima davao je dobre rezultate za klase koje nemaju premali broj primera.

Kada je u pitanju LSTM model za klasifikaciju, nad podacima iz trening skupa isproban je SMOTE (eng. *Synthetic Minority Oversampling TEchnique*) algoritam za balansiranje podataka. Ova tehnika podrazumeva stvaranje novih sintetičkih primera kod manjinskih klasa. Zasniva se na *k-nearest neighbors* metodi, gde se najpre bira jedan slučajni primer iz klase, a zatim i njegovih *k* najbližih suseda. Potom se od pronađenih suseda, na slučajan način bira jedan podatak, nakon čega se stvara sintetički primer koji odgovara slučajno odabranoj tački između 2 izabrana podatka u prostoru obeležja[10]. Nakon primene ove tehnike za balansiranje, svaka od klasa ima 1783 primera (jednako broju primera prethodno najbrojnije klase u trening skupu – *Parties*).

Tabela 1. *f1 micro score* nad testnim skupom

Model	F1 micro score
BERT bez <i>oversampling-a</i>	0.85
BERT + prvi način <i>oversampling-a</i>	0.89
BERT + drugi način <i>oversampling-a</i>	0.86
LSTM bez <i>oversampling-a</i>	0.72
LSTM + SMOTE <i>oversampling</i>	0.51

U tabeli 1 su prikazane ostvarene performanse isprobanih modela dobijene nad testnim skupom podataka.

Mera evaluacije koje je korištena za sve navedene modele, kao i za njihovo međusobno poređenje je *f1_micro* mera. Ona se najčešće koristi za proveru kvaliteta *multi-label* binarnih problema, što je i ovde slučaj. Računa se kao srednja vrednost *f1_score* mere svih klasa. *F1_score* se računa po formuli:

$$2 \times ((Precision \times Recall) \div (Precision + Recall)).$$

Precision predstavlja odnos tačno pozitivnih predikcija (eng. *True positives*) i svih pozitivnih predikcija. Ako je broj netačno pozitivnih predikcija veći od broja tačno pozitivnih predikcija model će imati manji *precision*. Sa druge strane *Recall* predstavlja meru tačno pozitivnih predikcija u odnosu na sve tačno pozitivne primere. Idealna *f1 micro* mera iznosi 1, što bi značilo da model savršeno klasifikuje klauzule.

VI. ZAKLJUČAK

U ovom radu razmatrali smo i isprobali različita rešenja za problem *multi-label* klasifikacije klauzula iz pravnih ugovora. Motivacija za izradu ovog sistema predstavlja to što su komercijalni ugovori često veoma dugi i kompleksni tekstovi, te bi automatizacija procesa praćenja i pregleda ugovora bila od velikog značaja za pravnike.

Skup podataka čini preko 12000 labeliranih pravnih klauzula u jednu ili više kategorija. Ukupan broj kategorija je 41. Nad podacima smo najpre primenili pretprocesiranje i *word embedding*. Za *word embedding* smo isprobali dve različite tehnike: GloVe i BERT. Nakon toga smo za postupak klasifikacije klauzula trenirali i testirali više modela: BERT transformer model sa i bez balansiranja podataka i LSTM rekurentnu neuronsku mrežu sa i bez balansiranja podataka. Model smo evaluirali nad testnim skupom. Mera evaluacije koja je korištena je *f1 micro* mera. Najbolji rezultat je ostvario BERT model uz *oversampling* podataka (0.89 i 0.86). Primetno je da je LSTM model dao znatno slabije rezultate

nakon primene SMOTE algoritma za balansiranje podataka, na osnovu čega možemo zaključiti da sintetički primeri nisu dovoljno relevantni u ovom slučaju. Razlog tome je verovatno izuzetno mali broj podataka koji poseduju manjinske klase. Sa druge strane, *RandomOverSampler* se pokazao kao bolje rešenje, jer su se rezultati koje je ostvario BERT model poboljšali.

Potencijalno poboljšanje performansi u odnosu na one koje postižu isprobani modeli bi se moglo postići najpre proširenjem skupa podataka: dodavanjem novih pravnih ugovora i klauzula. Osim toga, primena ansambl binarnih modela (transformer i/ili RNN) bi potencijalno mogla ostvariti bolje performanse. U tom slučaju bi se trenirao 41 model za svaku od kategorija. Razlog zbog koga ovaj pristup nije primenjen su hardverska ograničenja. Takođe, kompleksnije arhitekture neuronskih mreža uz optimizaciju više parametara bi isto tako mogle ostvariti bolje performanse.

REFERENCES

- [1] Dale, Robert. (2019). Law and word order: NLP in legal tech. *Natural Language Engineering*. 25. 211-217. 10.1017/S1351324918000475.
- [2] How AI Is Being Used In The Legal Industry, *Forbes*, poslednji pristup: 30.03.2021.,
url:
<https://www.forbes.com/sites/forbesbusinesscouncil/2021/01/19/how-ai-is-being-used-in-the-legal-industry/?sh=634b6e1550c6>
- [3] Hendrycks, Dan & Burns, Collin & Chen, Anya & Ball, Spencer. (2021). CUAD: An Expert-Annotated NLP Dataset for Legal Contract Review.
- [4] Chalkidis, Ilias & Fergadiotis, Manos & Malakasiotis, Prodromos & Aletas, Nikolaos & Androutsopoulos, Ion. (2020). LEGAL-BERT: The Muppets straight out of Law School. 2898-2904. 10.18653/v1/2020.findings-emnlp.261.
- [5] G. Endel, F. and Piringer, H. (2015) 'Data Wrangling: Making data useful again', *IFAC-PapersOnLine*, 48(1), pp. 111–112.
- [6] Ms. Anjali Ganesh Jivani, A Comparative Study of Stemming Algorithms, Anjali Ganesh Jivani et al, *Int. J. Comp. Tech. Appl.*, Vol 2 (6), 1930-1938, ISSN:2229-6093.
- [7] Raghavan A K, Venkatesh Umaashankar, and Gautham Krishna Gudur. 2019. Label frequency transformation for multi-label multi-class text classification. In *GermEval 2019, 15th Conference on Natural Language Processing (KONVENS 2019)*, Erlangen, Germany. R. Nicole, "Title of paper with only first word capitalized," *J. Name Stand. Abbrev.*, in press.
- [8] Tsoumakas G, Katakis I. Multi-label classification: An overview. *Dept of Informatics, Aristotle University of Thessaloniki; Greece: 2006*. M. Young, *The Technical Writer's Handbook*. Mill Valley, CA: University Science, 1989.
- [9] Xiaofei Ma, Zhiguo Wang, Patrick Ng, RameshNallapati, and Bing Xiang. 2019. Universal Text Representation from BERT: An Empirical Study. *arXiv:1910.07973 [cs]*.
- [10] Chawla, N., Bowyer, K., Hall, L., & Kegelmeyer, W. SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 2002, 16: 341-378.