

Analiza emocija rekurentnim neuronskim mrežama

Seminarski rad u okviru kursa
Metodologija stručnog i naučnog rada
Matematički fakultet

Marko Vićentijević, Sanja Mijatović, Sanja Radulović
markovicentijevic2702@gmail.com, mijatovicsanja000@gmail.com,
sanjaradulovic997@gmail.com

12. decembar 2021.

Sažetak

Emocija je ideja, mišljenje ili osećaj o nekoj temi. Zadatak analize emocija je odrediti koja emocija se izražava u zadatom tekstu. Najčešće se koncipira kao problem klasifikacije tekstova na pozitivne i negativne. Rad prikazuje nekoliko varijanti rekurentnih neuronskih mreža primenjenih na problem klasifikacije filmskih recenzija domaćih filmova koje su predstavljene kratkim rečenicama na srpskom jeziku u pozitivne i negativne. Prikazan je čitav proces: prečišćavanje podataka, primeri koda za arhitekturu svakog od modela, analiza rezultata i poređenje predstavljenih modela.

Sadržaj

1	Uvod	3
2	Podaci	3
3	Pretprocesiranje podataka	4
3.1	Učitavanje podataka	4
3.2	Lematizacija	4
3.3	Uklanjanje stop reči	5
3.4	Tokenizacija	5
4	Rekurentna neuronska mreža	6
4.1	Potpuno povezana rekurentna neuronska mreža	6
4.2	Dugo kratkoročna memorija	7
4.3	Rekurentna jedinica sa kapijama	8
4.4	Bidirekciona asocijativna memorija	8
5	Rezultati	9
6	Zaključak	12
	Literatura	13

Listings

1	Potpuno povezana rekurentna neuronska mreža	6
2	Jednoslojna LSTM rekurentna neuronska mreža	7
3	Troslojna LSTM rekurentna neuronska mreža	8
4	GRU rekurentna neuronska mreža	8
5	LSTM Bidirekciona rekurentna neuronska mreža	8
6	GRU Bidirekciona rekurentna neuronska mreža	9

1 Uvod

Obrada prirodnih jezika (eng. Natural language processing - NLP) predstavlja podoblast lingvistike, računarstva i veštačke inteligencije koja se bavi interakcijama između računara i ljudskog jezika. Cilj ove metode je da računar na neki način razume sadržaj dokumenta uključujući kontekstualne nijanse jezika u dokumentu. Tehnologija tada može precizno izdvojiti informacije i uvide sadržane u dokumentima, kao i kategorizaciju i organizaciju samog dokumenta. Prirodni jezici nemaju inherentna ograničenja za razliku od programskih jezika koji su unapred strogo definisani, takođe tekstovi na prirodnom jeziku odlikuju se visokom kompleksnošću, po mogućstvu mogu biti dvosmisleni što dodatno daje na težini obrade dokumenta prirodnog jezika od strane računara.

Metode zasnovane na ručno sastavljenim pravilima i bazama znanja (eng. Rule-based NLP) (1950 - 1990) predstavlja sam početak razvijanja NLP metoda. Skupovi pravila kreirani su od strane stručnjaka iz proučavane oblasti i predstavlja formu pravila koje bi računar pratio. Pristup se pokazao kao veoma spor i zahtevan.

Metode zasnovane na statistici i podacima (eng. Statistical NLP) (1990 -) počele su da se koriste od 1990. kada se razvijanje NLP-a usmerilo statističkom pristupu zasnovanom na nadgledanom mašinskom učenju, gde nisu eksplicitno navođene veze između podataka i predikcija, već se baziralo na parovima rečenica i tačnih izlaza.

Neuronski NLP (eng. Neural NLP) počeo je da se primenjuje od 2010. godine, kada primena neuronskih mreža u obradi prirodnog jezika uzima široku primenu. NLP doživljava ovu transformaciju delimično zbog toga što ove tehnike daju najpreciznije rezultate.

Analiza emocija (eng. Sentiment analysis) ili istraživanje mišljenja predstavlja proces računarske identifikacije i kategorizacije mišljenja izraženih u delu teksta, kako bi se utvrdilo da li je stav pisca prema određenoj temi pozitivan, negativan ili neutralan. To je jedna od najaktivnijih istraživačkih oblasti u obradi prirodnog jezika i pretrazi teksta poslednjih godina [8]. Analiza emocija je upotreba NLP-a, računarske lingvistike i biometrije za sistemsko identifikovanje, izdvajanje, kvantifikovanje i proučavanje afektivnih stanja i subjektivnih informacija.

U ovom radu bazirali smo se na analizi emocija filmskih recenzija na srpskom jeziku koristeći rekurentne neuronske mreže. Cilj rada je dobiti matematički model koji sa visokim procentom tačnosti određuje da li je filmska recenzija napisana u pozitivnom ili negativnom kontekstu.

2 Podaci

Podaci su preuzeti sa [SentiComments.SR](#) koji ujedno i predstavlja skup izvučenih filmskih recenzija sa domaćih filmskih sajtova na srpskom jeziku. Podaci su ručno klasifikovani sa **čtetvorovalentnom notacijom** koja predstavlja **objektivnost, subjektivnost, sarkazam i polarnosti**. U ovom radu podaci su modifikovani prema **polarnoj notaciji** u kojoj nam je bitno samo da li je recenzija pozitivna ili negativna ne osvrćemo se na subjektivnost, negativnost ili sarkazam odnosno dodeljen im je **pozitivan**

kontekst (1) ili **negativan kontekst (0)**. Ovaj skup podataka se sastoji od 3490 filmskih recenzija. Deo učitanih podataka je prikazan na slici 1. Kolona **sentiment_label** sadrži oznake osećanja +1, -1, +M, -M, +NS, -NS. U prikazanim podacima imamo prisustvo oznake +1 koja se koristi za filmove koji su potpuno pozitivni i oznake +NS koja se koristi za tekstove koji su objektivni, ali u binarnoj klasifikaciji (pozitivno, negativno) više naginju na pozitivno. Kolona **comment_id** predstavlja jedinstveni identifikator na osnovu filma i njegovog predstavljanja u drvetu komentara za taj određeni film i time se označava njegova pozicija u čitavoj diskusiji. Kolona **comment_text** predstavlja tekstove komentara.

	sentiment_label	comment_ID	comment_text
0	+1	1-1	♥ Znao sam da će ovaj biti prvi! :)
1	+NS	1-2	pa mora... The Dude Abides! :)
2	+1	1-6	Film gledam već godinama i svaki put otkrijem ...
3	+1	1-7	Svaki faktor je podjednako krucijalan i savrše...
4	+1	1-8	john goodman je obeležio ovaj film sa svojim p...
5	+1	1-9	Hvala za info da je večeras bio film. Odličan ...
6	+1	1-11	Sjajan odabir za prvi tekst. :) Odličan film, ...
7	+1	1-12	verovatno najsmešniji film svih vremena :)
8	+1	1-13	sjajan film
9	+1	3-1	10/10, oduševljen sam filmom. Sve što si napis...

Slika 1: Učitani podaci

3 Pretprocesiranje podataka

Neuronska mreža na ulazu ne može da primi matricu rečenica, već na ulazu očekuje vektor brojeva za svaki komentar koji ulazi u mrežu i zato je potrebno određeno pretprocesiranje od niza filmskih recenzija do matrice numeričkih vektora koji predstavljaju filmske recenzije. U ovoj glavi predstavimo detaljno u fazama pretprocesiranje naših podataka.

3.1 Učitavanje podataka

U učitanoj skupi podataka **SentiComments.SR** zbog konzistentnosti sve filmske recenzije su prevedene sa ćirilice na latinicu. Tabela podaci su podeljeni na filmske recenzije i njihovo emocionalno značenje koje je obeleženo četvorovalentnom notacijom. Potom četvorovalentna notacija je prebačena u polarnu notaciju koja je predstavljena sa **1 (pozitivan komentar)** i **0 (negativan komentar)**. U filmskim recenzijama emotikoni (eng. emoticon) se obrađuju tako da budu konzistentni i interpunkcijski znakovi su uklonjeni.

3.2 Lematizacija

Zbog ograničenosti količine podataka reči se svode na njihov kanonski oblik, čime se bavi **lematizacija**. Ovim je omogućeno da mreža ima više

reči u različitim rečenicama. Istraživanje nas je dovelo do [classla](#) Python paketa koji vrši lematizaciju reči standardnog i nestandardnog srpskog jezika sa tačnošću lematizacije od približno 99%.

3.3 Uklanjanje stop reči

U pozitivnim i u negativnim komentarima nalaze se iste reči koje filmskim recenzijama ne daju ni pozitivan, a ni negativan kontekst, takve reči nazivaju se **stop reči**. Jedna od metoda za uklanjanje stop reči je metoda TF-IDF. **TF-IDF (Term Frequency - Inverse Document Frequency)** je pristup ponderisanja termina koji se obično koristi za predstavljanje tekstualnih dokumenata kao vektora [9]. TF-IDF ima za cilj koliko je reč važna u dokumentu, na primer za jednu filmsku recenziju u skupu filmskih recenzija. U našem projektu pronašli smo [skup stop reči](#).

3.4 Tokenizacija

Nakon postupka lematizacije i uklanjanja stop reči, potrebna je još podela teksta na tokene (reči i ostale smislene elemente), odnosno uraditi pretvaranje niza reči u vektore celih brojeva. U ovoj fazi može se koristiti ugrađena klasa iz biblioteke `tensorflow` za pretprocesiranje podataka koja se naziva `Tokenizer`. `Tokenizer` primenjuje dve metode `fit_on_texts` i `texts_to_sequences`.

`Fit_on_texts` kreira vokabular baziran na frekvenciji reči.

`Texts_to_sequences` transformiše svaku reč iz teksta u sekvencu celog broja.

Ovo je poslednja faza nakon koje su naši podaci spremni za obradu od strane rekurentnih neuronskih mreža (eng. Recurrent Neural Network - RNN). Primer podataka nakon tokenizacije prikazan je na slici 2.

```
Encoded X Train
[[ 369   1 131 ...   28   54   31]
 [  18   1  43 ...   22   1 2644]
 [   1  33  63 ...    0    0    0]
 ...
 [ 268 514   0 ...    0    0    0]
 [  35   4  16 ... 6296   5 759]
 [  12   6 899 ... 643 6298   3]]

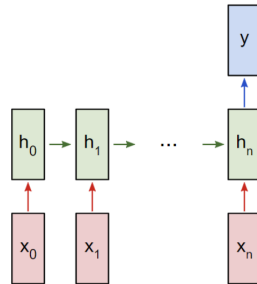
Encoded X Test
[[   5   4 164 ...    0    0    0]
 [ 531 1012  34 ...   17 5110   44]
 [  12  291 4667 ...    0    0    0]
 ...
 [  72 699   3 ...  131    0    0]
 [ 318 3842   3 ...    0    0    0]
 [  20   1   0 ...    0    0    0]]

Maximum review length: 16
```

Slika 2: Tokenizovani podaci

4 Rekurentna neuronska mreža

Rekurentna neuronska mreža predstavlja klasu veštačkih neuronskih mreža (eng. Artificial neural networks - ANN). Čvorovi ove neuronske mreže mogu se predstaviti usmerenim grafom duž vremenskog perioda. Ovo svojstvo omogućava mreži da pokaže svojstvo vremenskog dinamičkog ponašanja. Ove mreže mogu koristiti svoje unutrašnje zapamćeno stanje (memoriju) za obradu sekvenci ulaza promenljive dužine. Zato su često korišćene u zadacima povezanog prepoznavanja rukopisa i prepoznavanja govora. Izraz rekurentna neuronska mreža koristi se neselektivno za označavanje dve široke klase mreža sa sličnom opštom strukturom. Obe klase mreža pokazuju vremensko dinamično ponašanje. Rekurentna mreža sa konačnim impulsom je usmereni aciklični graf koji se može zameniti strogo naprednom neuronskom mrežom, dok se mreža sa beskonačnim impulsom koja je predstavljena usmerenim cikličnim grafom ne može zameniti. Obe ove mreže i sa konačnim impulsom i sa beskonačnim mogu imati unutrašnja (memorijska) stanja. Unutrašnja stanja se takođe mogu zameniti sa drugom mrežom ako ta mreža sadrži vremenska odlaganja ili ima povratne veze. Ovakva kontrolisana stanja nazivaju se zatvorena stanja ili zatvorena memorija i deo su **Long Short Term Memory - LSTM**. Primer opšte rekurentne neuronske mreže je prikazan na slici 3.



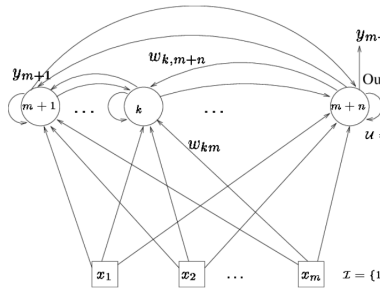
Slika 3: Slikovit prikaz RNN mreže

4.1 Potpuno povezana rekurentna neuronska mreža

Potpuno povezana rekurentna neuronska mreža (eng. Fully recurrent neural networks - FRNN) povezuje izlaze svih neurona sa ulazima svih neurona. Ovo predstavlja najopštiju topologiju neuronskih mreža jer se sve ostale topologije mogu prikazati kao potpuno povezujuće mreže tako što ćemo određenim težinama koje ne postoje dodeliti težinu nula. FRNN je prikazan na slici 4.

```
1000 model = Sequential()
      model.add(Embedding(input_dim = total_words, output_dim =
                           EMBED_DIM))
1002 model.add(SimpleRNN(EMBED_DIM, return_sequences = True))
      model.add(SimpleRNN(EMBED_DIM))
1004 model.add(Dense(1))
```

Listing 1: Potpuno povezana rekurentna neuronska mreža

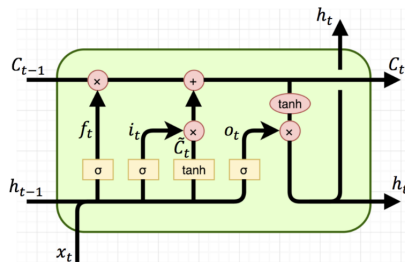


Slika 4: Slikovit prikaz FRNN mreže

4.2 Dugo kratkoročna memorija

Dugo kratkoročna memorija (eng. Long Short Term Memory - LSTM) je jedna od arhitektura rekurentnih neuronskih mreža, koja se koristi u dubokom učenju (eng. Deep Learning). Za razliku od standardnih naprednih neuronskih mreža, LSTM ima povratne veze. Pored pojedinačnih podataka (slike) LSTM može da obradi čitave sekvence podataka (glasovni poziv, video snimak). Primenljiv na zadatke kao što su povezano prepoznavanje rukopisa, prepoznavanje govora, otkrivanje anomalija u mrežnom saobraćaju. LSTM je prikazana na slici 5.

Uobičajena LSTM jedinica sastoji se od ćelije, ulazne kapije (eng. input gate), izlazne kapije (eng. output gate) i kapije za zaborav (eng. forget gate). Ćelija pamti vrednosti u proizvoljnim vremenskim intervalima, a tri kapije (eng. gates) regulišu protok informacija u ćeliju i iz nje. LSTM mreže su zasnovane na podacima iz vremenskih serija, jer može doći do zaostajanja nepoznatog trajanja između važnih događaja u vremenskoj seriji. LSTM je razvijena za rešavanje problema nestajućeg gradijenta na koji se može naići pri obuci tradicionalnih rekurentnih neuronskih mreža.



Slika 5: Slikovit prikaz LSTM mreže

```

1000 model = Sequential()
      model.add(Embedding(total_words, EMBED_DIM, input_length =
                           max_length))
1002 model.add(LSTM(LSTM_OUT))
      model.add(Dense(1, activation='sigmoid'))

```

Listing 2: Jednoslojna LSTM rekurentna neuronska mreža

```

1000 model = Sequential()
      model.add(Embedding(total_words, EMBED_DIM, input_length =
                           max_length))
1002 model.add(LSTM(units = 50, return_sequences = True, activation = '
      relu'))
      model.add(Dropout(0.2))
1004 model.add(LSTM(units = 50, return_sequences = True, activation = '
      relu'))
      model.add(Dropout(0.2))
1006 model.add(LSTM(units = 50, activation = 'relu'))
      model.add(Dropout(0.2))
1008 model.add(Dense(1, activation='sigmoid'))

```

Listing 3: Troslojna LSTM rekurentna neuronska mreža

4.3 Rekurentna jedinica sa kapijama

Zbog svoje otpornosti na nestajući eksplodirajući gradijent LSTM je najpopularnija i najkorišćenija vrsta rekurentnog sloja danas. Kako se sastoji iz četiri jednosmerna sloja bilo je brojnih pokušaja da se pojednostavi, a da se sačuva originalna funkcionalnost. Najpopularniju novu opciju izneli su K.Cho i saradnici 2014. koja se naziva rekurentna jedinica sa kapijama (eng. Gated Recurrent Unit).

Kapije ulaza i zaborava iz LSTM ćelije su spojene u jednu kapiju koja se sada naziva kapija osveživanja (eng. update gate), pod pretpostavkom da se najbitnije promene stanja ćelije dešavaju na istim neuronima u sloju. Takođe su u jedan sloj spojeni stanje ćelije i izlaza.

```

1000 model = Sequential()
      model.add(Embedding(input_dim=total_words, output_dim=EMBED_DIM))
1002 model.add(GRU(256, return_sequences=True))
      model.add(SimpleRNN(128))
1004 model.add(Dense(1, activation = 'sigmoid'))

```

Listing 4: GRU rekurentna neuronska mreža

4.4 Bidirekciona asocijativna memorija

Bidirekciona asocijativna memorija (eng. Bidirectional associative memory - BAM) predstavljena od strane Bart Kosko [5], predstavlja varijantu Hopfildove mreže koja skladišti asocijativne podatke u obliku vektora. Hopfildova mreža predstavlja memoriju koja obično mapira ulazni vektor u memorisani vektor čije je heming (eng. hamming) rastojanje od ulaznog vektora najmanje. Dvosmernost dolazi prolaskom informacija kroz matricu i njihovom transpozicijom. BAM mreža ima dva sloja, od kojih se svaki može pokrenuti kao ulaz za opoziv asocijacije i proizvesti izlaz na drugom sloju.[3]

```

1000 model = Sequential()
      model.add(Embedding(total_words, EMBED_DIM))
1002 model.add(Bidirectional(LSTM(LSTM_OUT)))
      model.add(Dense(64, activation = 'relu'))
1004 model.add(Dense(1, activation = 'sigmoid'))

```

Listing 5: LSTM Bidirekciona rekurentna neuronska mreža


```

1000 model = Sequential()
      model.add(Embedding(total_words, EMBED_DIM))
1002 model.add(Bidirectional(GRU(256)))
      model.add(Dense(64, activation = 'relu'))
1004 model.add(Dense(1, activation = 'sigmoid'))

```

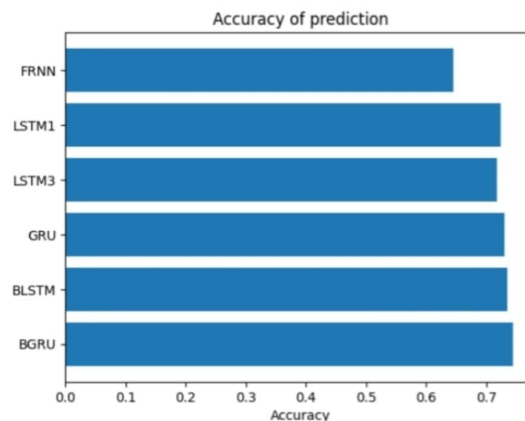
Listing 6: GRU Bidirekciona rekurentna neuronska mreža

5 Rezultati

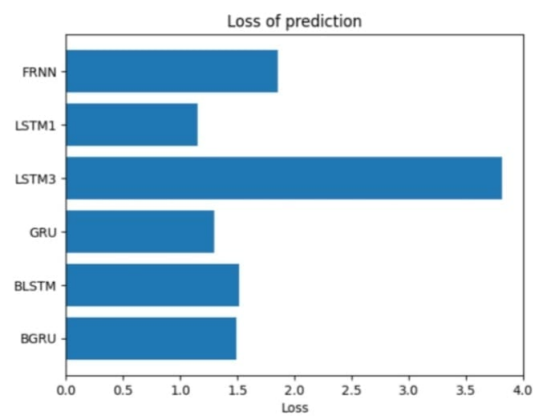
Kreirali smo šest modela rekurentnih neuronskih mreža čije su strukture prethodno gore prikazane.

1. **FRNN** - Potpuno povezanu rekurentnu neuronsku mrežu.
2. **LSTM1** - Jednoslojnu LSTM rekurentnu neuronsku mrežu.
3. **LSTM3** - Troslojnu LSTM rekurentnu neuronsku mrežu.
4. **GRU** - GRU rekurentnu neuronsku mrežu.
5. **BLSTM** - Bidirekcionu LSTM rekurentnu neuronsku mrežu.
6. **BGRU** - Bidirekcionu GRU rekurentnu neuronsku mrežu.

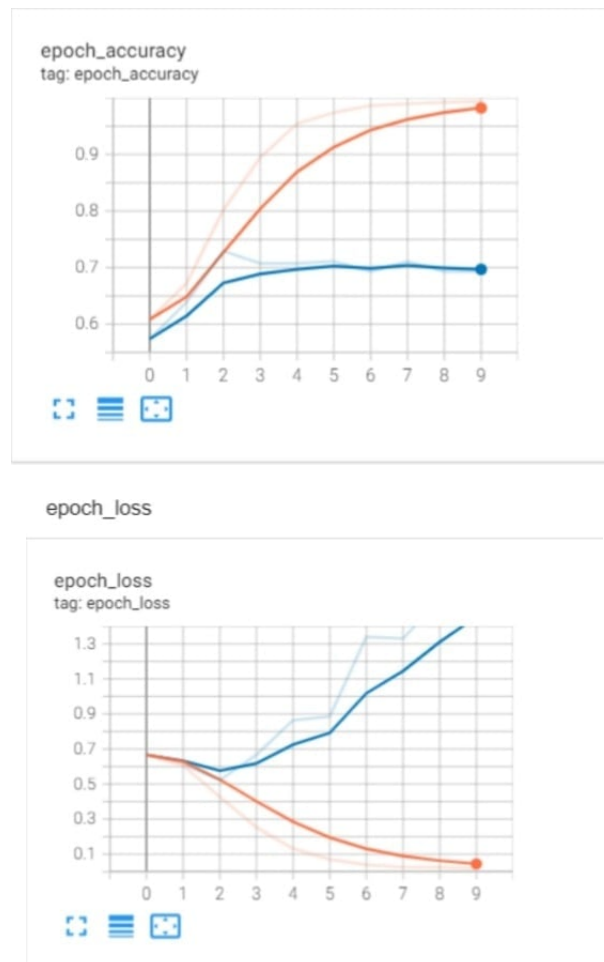
Kao što možemo videti po graficima koji su prikazani na slici 6 i na slici 7, najbolje se pokazala **BGRU** - Bidirekciona GRU rekurentna neuronska mreža sa **tačnošću od 0.7442** i **vrednosti funkcije greške od 1.4928**. Na slici 8 prikazane su vrednosti funkcije greške i tačnost po epohama.



Slika 6: Preciznost pri predikcijama



Slika 7: Funkcija greške pri predikcijama



Slika 8: Grafik koji prikazuje preciznost i funkciju greške BGRU modela kroz epohe

6 Zaključak

Ovim radom smo predstavili neke osnovne metode i koncepte koji se koriste za analizu emocija preko rekurentnih neuronskih mreža. Kombinovanjem eksperimentalnog pristupa i prakse zasnovane na naučnoj osnovi, identifikovane su tehnike pretprocesiranja i izgradnje raznih modela rekurentnih neuronskih mreža, sa ciljem dobijanja što veće preciznosti pri klasifikaciji filmskih recenzija. Projekat može biti unapređen većom količinom ulaznih podataka koji bi se koristili za treniranje mreže kao i pribegavanju određenim neuronskim mrežama koje se bave pretprocesiranjem podataka.

Literatura

- [1] Vuk Batanović. Metodologija rešavanja semantičkih problema u obradi kratkih tekstova napisanih na prirodnim jezicima sa ograničenim resursima, 2020. online na: <https://vukbatanovic.github.io/pdf/Disertacija.pdf>.
- [2] Jason Champion. Serbian stop words, 2016.
- [3] B. Coppin. *Artificial Intelligence Illuminated*. Jones and Bartlett illuminated series. Jones and Bartlett Publishers, 2004.
- [4] Christopher C.R. Turk David M.W. Powers. *Machine Learning of Natural Language*. Springer, London, 1989.
- [5] B. Kosko. Bidirectional associative memories. *IEEE Transactions on Systems, Man, and Cybernetics*, 18(1):49–60, 1988.
- [6] Bates M. Models of natural language understanding. 1995.
- [7] Kaja Dobrovoljc Nikola Ljubešić. *What does Neural Bring? Analysing Improvements in Morphosyntactic Annotation and Lemmatization of Slovenian, Croatian and Serbian*, volume Proceedings of the 7th Workshop on Balto-Slavic Natural Language Processing. Association for Computational Linguistics, 2019.
- [8] Ajit kumar P.G. Preethi, V. Uma. Temporal sentiment analysis and causal rules extraction from tweets for event prediction. *Procedia Computer Science*, 48:84–89, 2015. International Conference on Computer, Communication and Convergence (ICCC 2015).
- [9] Claude Sammut and Geoffrey I. Webb, editors. *TF-IDF*, pages 986–987. Springer US, Boston, MA, 2010.
- [10] Vahid Mirjalili Sebastian Raschka. *Python Machine Learning: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow 2, 3rd Edition*. Packt Publishing, 2019.
- [11] Boško Nikolić Vuk Batanović, Milos Cvetanović. A versatile framework for resource-limited sentiment articulation, annotation, and analysis of short texts. 2020.