



PREDIKCIJA VREDNOSTI KRIPTOVALUTA ANALIZOM ISTORIJSKIH CENA, BLOKČEJN INFORMACIJA I SENTIMENTA TVITOVA

PREDICTING THE VALUES OF CRYPTOCURRENCIES USING HISTORICAL VALUES, BLOCKCHAIN INFORMATION AND TWITTER SENTIMENT

Milica Milutinović, *Fakultet tehničkih nauka, Novi Sad*

Oblast – ELEKTROTEHNIKA I RAČUNARSTVO

Kratak sadržaj – *Ovaj rad bavi se problemom predikcije cene Bitkoina na osnovu istorijskih cena i blokčejn informacija. Cenu kriptovaluta velikim delom diktiraju špekulacije pa je ispitano u kojoj meri sentiment analiza tvitova doprinosi padu ili rastu cene. Pored podataka o ceni i sentimentu, korišćeni su blokčejn podaci i istorijski podaci za još tri popularne kriptovalute: Lajtkoin, Ethereum i Ripple. Za sentiment analizu ispitane su tri tehnike: konvolucione neuronske mreže, ansambl leksikona i metoda zasnovana na jezičkim pravilima. Rezultati sentiment analize upotrebljeni su u daljem procesu predviđanja kriptovaluta. Izabrani model za predikciju je rekurentna neuronska mreža sa GRU ćelijama. U fazi evaluacije posmatrala se tačnost predviđanja kretanja cene, odnosno da li će ona da raste ili opada. Takođe, upotrebljena je i relativna tačnost u kojoj se posmatra odnos stvarne cene i prediktovane. Najbolji rezultat dostigao je 57,3% tačnosti u predikciji kretanja cene, odnosno 99,39% relativne tačnosti.*

Ključne reči: sentiment analiza; predikcija; Bitkoin;

Abstract – *This paper studies the problem of Bitcoin price prediction based on historical prices and blockchain information. As the price of a cryptocurrency is dictated largely by speculation we also examine the influence of twitter sentiment on the performance of our predictive models. We have also experimented with historical data for three other popular cryptocurrencies: Litecoin, Ethereum and Ripple. We experimented with three techniques for sentiment mining: convolutional neural networks, ensemble of lexicons and linguistic rules. The results of sentiment analysis were integrated with other data to train a recurrent neural network with GRU cell as our price prediction model. We evaluated our model both as a binary classifier (predicting whether the price will go up or down) and a regression predictor (predicting the actual price). The accuracy of the classification model was 57.3%, while the relative accuracy of the regression model was 99.39%.*

Keywords: sentiment analysis; prediction; Bitcoin;

1. UVOD

Kriptovalute predstavljaju oblik digitalne imovine koje su zasnovane na kriptografskom protokolu koji omogućavaju korisnicima da ih čuvaju i razmenjuju putem mreže [1]. Satoshi Nakamoto je u januaru 2009. implementirao Bitkoin (en. *Bitcoin*), što predstavlja prvu kriptovalutu. Od tada njihova popularnost stalno raste, najviše zbog osobina kao što su decentralizovanost, transparentnost i sigurnost. U ovom radu će biti opisani pristupi i rešenja za predviđanje cene Bitkoina. U poređenju sa „klasičnim“ akcijama na berzi, za kriptovalute se mogu vezati informacije o samoj blokčejn tehnologiji koje su javno dostupne što uvodi dodatnu semantiku [2]. U narednoj sekciji je izložen pregled postojeće literature vezane za problem predviđanja cena kriptovaluta i sentiment analize podataka. Treća sekcija sadrži informacije o prikupljanju podataka. Četvrta sekcija opisuje metodologije i alate korišćene u radu. U petoj sekciji objašnjena je eksperimentalna evaluacija rezultata, dok poslednja sekcija zaključuje rad predlaže pravce mogućeg daljeg istraživanja.

2. PRETHODNA REŠENJA

Dosadašnja literatura pokriva veliki broj problema koji teže da se objasne u radu. Proučavanja sentiment analize teksta su dospila visok nivo, a pristupi koji su rađeni variraju od metoda koje ne koriste nikakav vid učenja odnosno koja se baziraju na pravilima i/ili leksikonima, do metoda koje koriste mašinsko učenje za te svrhe. Metod zasnovan na pravilima [3], VADER alat, teži da odredi sentiment na osnovu jezičkih konstrukcija, odnosno pravila koje vladaju u njima. Zbog izuzetnih rezultata koji su dobijeni (0,96 F1 mera tačnosti za sentiment tvitova), implementacija ovog metoda je korišćena kao jedan od načina za određivanje sentimenta. Pristup koji određuje sentiment teksta baziran na ansamblima leksikonima [4] traži reči u formiranim leksikonima i u zavisnosti od broja ponavljanja u njima, dodeljuju svakoj reči jedinstven sentiment. Ovaj metod unapredio je proste metode bazirane na leksikonima koji su imali mnogih nedostataka [5], a dospao je približnu tačnost kao i metode mašinskog učenja u određivanju sentimenta teksta. Iz tih razloga, implementiran je, a detalji će biti objašnjeni kasnije u radu. Analiziranje sentimenta u tekstu često se rešava primenom algoritama dubokog učenja. Rad zasnovan na konvolucionim neuronskim mrežama [6] koristi konvolucionu arhitekturu neuronske mreže i tzv. Bag-of-Words model [7] kao

NAPOMENA:

Ovaj rad proistekao je iz master rada čiji mentor je bio dr Aleksandar Kovačević, van. prof.

ulazne podatke u mrežu. Iako arhitektura u pomenutom radu nije velike dubine, podešavanjem parametara mreže i korišćenjem regularizacije u procesu učenja, ovaj model pokazao je bolje rezultate od drugih koji su upoređivani u istom radu.

U radu koji se bavi predikcijom volatilnosti akcija na kineskom tržištu [8], korišćenjem rekurentne neuronske mreže sa sentimentalnim podacima, dobijena je tačnost od 65%. Međutim, u tom radu predikcija je vršena samo za jedan vremenski korak unapred.

3. SAKUPLJANJE PODATAKA

Skup podataka koji je analiziran dobijen je prikupljanjem tvitova influensera u domenu kriptovaluta, istorijskih cena i informacija o blokčejnu u periodu od Januara 2016. godine do kraja Juna 2018. godine.

Tvitovi su prikupljeni zahvaljujući projektu [9] koji formira url na osnovu koga skrejpuje (eng. *scrapping*) stranicu preko koje sakuplja tvitove. Zahvaljujući ovom pristupu preuzet je veliki broj tvitova u kratkom vremenskom intervalu. Nakon procesa sentiment analize, dobijeni rezultati prosleđeni su u fazu predikcije. Iskorišćeno je 7 numeričkih atributa u daljem toku:

- vader_sent i vader_pol – intezitet i polaritet sentimenta dobijen korišćenjem VADER alata
- lex_sent - sentiment dobijen korišćenjem modela ansambla leksikona
- cnn_sent i cnn_pol - intezitet i polaritet sentimenta dobijen korišćenjem CNN arhitekture
- favorites - predstavlja broj "lajkovanja" tvita
- retweets – predstavlja broj deljenja tvita

Istorijski podaci o cenama su preuzeti sa Poloniex API-ja za Bitkoin, Lajtkoin, Ethereum i Ripl. Podaci se nalaze u vremenskom koraku od 5 minuta za period od Januara 2016. godine do Aprila 2018. godine. Eksplorativnom analizom atributa je utvrđeno da podaci o ceni imaju veliku međusobnu korelaciju (približno 1) te je za dalju obradu korišćeno: close (cena na kraju intervala), date (vreme u *unix timestamp*-ovima), volume (trgovani iznos u poslednja 24h u dolarima) i quoteVolume (trgovani iznos u poslednja 24h u datoru kriptovaluti). Takođe je zaključeno da postoji korelacija između različitih valuta, tako da su za konačnu predikciju cene Bitkoina validirani i podaci od drugih valuta.

Preostalih 10 atributa su podaci o Bitcoin blokčejnu koji su preuzeti sa blockchain.info API-ja i unose dodatnu semantiku o transakcijama ove kriptovalute.

S obzirom da prikupljeni tvitovi nemaju pravilan vremenski interval, a podaci o blokčejnu se nalaze na uniformnom intervalu od 1 dan, nije ih bilo moguće spojiti sa podacima o cenama koje imaju uniforman interval od 5 minuta. Stoga su oni interpolirani na isti vremenski raspon. U radu je upotrebljena interpolacija na 1 sat i na 10 minuta.

Za proces formiranja rečnika u metodologiji koja koristi ansamble leksikona upotrebljeno je 5 skupova podataka koji sadrže anotirane tvitove: *Stanford Twitter Sentiment Test Set* [10], *Sentiment Strength Twitter Dataset* [11], skup od 9637 tvitova skinutih sa *Kaggle*-a, skup od 50800 tvitova skinutih sa *Kaggle*-a i *Senticnet* [12].

4. METODOLOGIJA

U cilju istraživanja sentimenta tvitova, isprobana su tri popularna rešenja. Rezultati metoda su upotrebljeni da se odredi u kojoj meri doprinose poboljšanju predikcije.

U sva tri pristupa, tvitovi su preprocesirani primenom nekoliko koraka: izbacivanjem emotikona, akronima i skraćenica, pretvaranjem velikih u mala slova, uklanjanjem *URL*-ova, praznih mesta, stop reči, znakova interpunkcije i specijalnih karaktera. Potom je primenjivana morfološka normalizacija reči. Skupovi podataka korišćeni u procesu formiranja leksikona preprocesirani su na sličan način, s tim da su im i ciljne labele svedene na isti interval, odnosno 1 za pozitivan, -1 za negativan, a 0 za neutralan tvit.

Prvi isprobani pristup koristi jezička pravila za računanje sentimenta [3]. U pomenutom pristupu posmatra se grupa reči koja predstavlja neku celinu, pa se njoj određuje sentiment, a ne samo pojedinačna reč. Najbolji rezultati se dobijaju ukoliko se radi analiza tzv. mikro-blogova, odnosno kratkih statusa. Implementacija ovog pristupa je javno dostupna i preuzeta je.

Drugi način bazira se na formiranju leksikona, koji se potom koriste da se pronađu reči koje imaju pozitivan ili negativan sentiment. Ovom metodologijom najpre su generisani leksikoni koristeći pomenute skupove. Leksikoni su kreirani tako što su se pronalazile reči koje su najzastupljenije u pozitivnim, odnosno negativnim recenzijama. Na ovaj način, napravljen je skup od 10 rečnika za navedenih 5 skupova, odnosno rečnici sa po 5 i 25 najpozitivnijih i najnegativnijih reči za svaki skup posebno. Formiran je i jedan rečnik koji je agregira sve skupove podataka zajedno i među njima traži skup od 25 najpozitivnijih i najnegativnijih reči. Preostala 4 rečnika formirana su na osnovu reči koje jednoznačno doprinose promeni inteziteta sentimenta i preuzeti su iz rada [13].

Sledeći korak predstavlja klasifikaciju sentimenta na pozitivan ili negativan na osnovu ansambla rečnika. Ovaj model uzima pojedinačne reči iz rečenice i predviđa njihov sentiment tako što traži da li se ta reč pojavljuje više kao pozitivna ili u negativna. Ukoliko preovladava pojavljivanje reči u leksikonima sa negativnim rečima, vrednost sentimenta reči će biti -1, odnosno 1 u suprotnom. Ukoliko se reč ne pojavljuje ni u jednom leksikonu, vrednost njenog sentimenta će biti 0. Nakon što se svakoj reči dodeli sentiment, računa se ukupna vrednost sentimenta na nivou teksta, odnosno tvita u našem slučaju. Ukupan sentiment za leksikon l i tvit t računa se po formuli:

$$s_l = \begin{cases} 1, \wedge if \sum (l, t) > 0 \\ -1, \wedge if \sum (l, t) < 0 \\ 0, \wedge otherwise \end{cases}$$

Konačan rezultat sentiment analize za jedan tvit dobija se usrednjavanjem rezultata svih leksikona. Ukoliko je izračunata vrednost veća od 0, uzima se da je tvit pozitivan, ukoliko je manja od 0, smatra se da je tvit negativan, a u suprotnom je tvit neutralan.

Treći pristup korišćen u ovom radu zasniva se na primeni konvolucionih neuronskih mreža, koje koriste konvolucioni operator da proizvode nove osobine. Ulaz u mrežu dobijen je formiranjem tzv. Bag-of-Words modela i predstavlja matricu fiksne dužine $N \times K$, gde N predstavlja

broj tвитова у датом интервалу, а K фиксну дужину рећеница на којима је извршено употпуњавање (*padding*) или одсечење (*truncate*) у односу на фиксну дужину. Излаз из мреже дaje две вредности, тј. однос раста или пада цене и ознаку да ли цена расте или опада.

Како се о твитовима чувaju информације о дану када су kreirani, излазни подаци у процесу тренирања se добијају тако што se računa однос цена tog дана i sutrašnjeg дана i ознаку da li cena raste ili opada. Na primer, ukoliko je tвит postavljen дана 10.12.2017. posmatra se cena tog дана i cena 11.12.2017, računa se količnik cena ta dva дана, a postavlja se i индикатор да ли cena raste ili opada. Zbog nemogućnosti да se odredi čiji je tвит važniji, svi tвитови objavljeni истог дана имају isti однос цена i isti индикатор koji говори o kretanju cene kriptovalute.

Za proces računanja predikcije korišćene su rekurentне neuronske mreže jer imaju sposobност да adresiraju vremensке податке користеći unutrašnju memoriju.

S obzirom da su svi atributi numeričки, odrđena je normalizacija u intervalu $[0,1]$ na свим attributima. Како bi se kreirala veća količina podataka upotrebljen je preklapajući obučavajući i trening skup. Ova redundantnost omogućena je korišćenjem „klizećeg“ vremenskog okvira за svaki skup ulaznih i izlaznih вредности.

Odabir arhitekture RNN mreže je добијен експериментално nakon što je testirano više različitih arhitektura. GRU koristi manje parametara od LSTM (eng. *Long short-term memory*), a performanse se neznatno razlikuju. Korišćene su GRU ћелије које sadrže 96 neurona (ћелија) u jednom sloju. Na izlazu je Fully-Connected слој sa 24 izlaza koji predstavljaju vremensке кораке предикције. Како bi se спречио overfitting, dodata je Dropout regularizacija. Testirani su модели sa једним i два unutrašnja скриена слоја. За функцију грешке одабрана је MSE (*mean squared error*), где је takođe korišćena Adamova оптимизација.

5. EKSPERIMENTI I REZULTATI

Za основну меру тачности је узет број успешио предвиђених кретања валута – пад или раст што може да се сведе на бинарну класификацију. Formula по којој се она израчунава:

$$Acc = \frac{\text{successful predictions}}{\text{total predictions}}$$

Relativna тачност, добија се računanjem односа стварне цене i предиковане. Осим поменутих метрика, korišćene su i srednja kvadratna i srednja apsolutna грешка, radi upoređivanja тачности rezultata.

У поступку evaluacije kreirano je više skupova kako bi se validirali подаци добијени u поступку sentiment analize. Урађени su i testovi radi испитивања параметара који doveđe do poboljšanja rezultata. Tim putem испитано од колико je značaja vremenski raspon u kojem su interpolirani подаци, број ulaznih vremenskih корака u процес обуčавања, као i информације o drugim vrednostima kriptovaluta.

Interesantno je primetiti da информације o sentimentu influensera доводе до boljih rezултата u одређеној meri. Uočено je da povećanje броја ulazних јединица, као i информације o vrednostima других kriptovaluta ne doprinose boljim rezултатима. Ipak, највеће побољшање постигнуто je u slučaju

smanjenja интервала interpolacije sa 1. сата на 10 минута. U Tabeli 1. prikazani su rezultati за најбољи validirani test, односно за za slučaj interpolacije на 10 минута, где se prediktuju вредности за наредна 2 сата, на основу вредности u prethodna 24 сата. Oznake skupova сastoje se od atributa koji ih сачинjavaju: **b** (блокчјен подаци), **h** (историјске подаци), **s** (sentiment подаци), **v** (sentiment подаци добијени vader методом), **c** (sentiment подаци добијени cnn методом) i **I** (sentiment подаци добијени ansamblima лексикона).

Tabela 1. *Rezultati*

	acc (%)	rel_acc (%)	mse	mae
bps	53.31	98,83	9180.78	79.47
bp	52.74	97.8	27721.35	151.15
bpv	47.36	96.1	79079.49	272.38
bpl	54.32	99.2	5056.30	54.41
bpc	52.76	97.22	40463.36	190.12
s	71.22	38.57	56364.57	100.56
ps	57.3	99.39	3731.67	40.63



Slika 1. *Grafik predikcije za skup ps.*

Na slici 1. dat je prikaz најбољег скупа атрибута ps, за date параметре. Може се уочити да су предиктовани rezultati prilično korelirani sa stvarnim.

6. ZAKLJUČAK

U ovom radu predložen je нови приступ предвиђања cena kombinacijom različitih skupova података u циљу preciznijih rezultata. Подаци o sentimentalnoj analizi uticajnih ljudi на Tвiteru dopринели су побољшању предикције, међутим, потребно je dodatno istraživanje kako bi se испитao степен успешиости. Ovim rezultatima je потврђена težina i kompleksnost problema predviđaња цене на berzi. Moguća обrazлоžења lošijih rezultata je relativno mali period скупа података (2 године i 6 meseci) i чинjenica da je u proteklih godinu дана забележен раст неких kriptovaluta čак i до 10,000% te je створен tržišni balon (*market bubble*) који je полако почео да puca u januaru 2018. године.

Budući првци istraživanja bi могли да размотре korišćenje неких eksperimentalnih модела neuronskih mreža u rešavanju problematike. Такође, било bi pogодно istražiti koliko pojedini tвитови утичу на промену цене, па u складу са tim da im se dodele odgovarajuće težine. Pored ovoga, било bi zanimljivo istražiti i како финансиске вести i неки други индикатори утичу i на друге цене на berzi i pored kriptovaluta.

7. LITERATURA

- [1] Wikipedia contributors. (2018, April 11). Cryptocurrency. In *Wikipedia, The Free Encyclopedia*. Retrieved 18:46, April 13, 2018, from <https://en.wikipedia.org/w/index.php?title=Cryptocurrency&oldid=835965914>
- [2] Blockchain. (2018, April 13). Retrieved April 13, 2018, from <https://en.wikipedia.org/wiki/Blockchain>
- [3] Hutto, C.J. & Gilbert, E.E. (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. Eighth International Conference on Weblogs and Social Media (ICWSM-14). Ann Arbor, MI, June 2014
- [4] Augustyniak, L., Kajdanowicz, T., Szymanski, P., Tuliglowicz, W., Kazienko, P., Alhajj, R., & Szymanski, B.K. (2014). Simpler is better? Lexicon-based ensemble sentiment classification beats supervised methods. 2014 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2014), 924-929.
- [5] S.-H. Na, Y. Lee, S.-H. Nam, and J.-H. Lee, "Improving opinion retrieval based on query-specific sentiment lexicon," in Advances in Information Retrieval, ser. Lecture Notes in Computer Science, vol. 5478. Springer Berlin / Heidelberg, 2009, pp. 734–738.
- [6] K., & Y. (2014, September 03). Convolutional Neural Networks for Sentence Classification. Retrieved April 13, 2018, from <https://arxiv.org/abs/1408.5882>
- [7] Wikipedia contributors. (2018, March 26). Bag-of-words model. In *Wikipedia, The Free Encyclopedia*. Retrieved 20:07, April 13, 2018, from https://en.wikipedia.org/w/index.php?title=Bag-of-words_model&oldid=832576977
- [8] Liu, Y., Qin, Z., Li, P., & Wan, T. (2017). Stock Volatility Prediction Using Recurrent Neural Networks with Sentiment Analysis. *Advances in Artificial Intelligence: From Theory to Practice Lecture Notes in Computer Science*, 192-201. Doi:10.1007/978-3-319-60042-0_22
- [9] J. (2018, January 29). Jefferson-Henrique/GetOldTweets-python. Retrieved April 13, 2018, from <https://github.com/Jefferson-Henrique/GetOldTweets-python>
- [10] Go, A., Bhayani, R., Huang, L.: Twitter sentiment classification using distant supervision. CS224N Project Report, Stanford (2009)
- [11] Thelwall, M., Buckley, K., Paltoglou, G.: Sentiment strength detection for the social web. *Journal of the American Society for Information Science and Technology* 63(1), 163–173 (2012)
- [12] SenticNet. (n.d.). Retrieved from <https://sentic.net/>
- [13] A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts, "Learning word vectors for sentiment analysis," in The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1. Association for Computational Linguistics, 2011, pp. 142–150.

Kratka biografija:



Milica Milutinović rođena je 6.1.1995. godine u Rumi. Osnovne studije je upisala 2013. godine na Fakultet tehničkih nauka, odsek Računarstvo i automatika. Osnovne studije je završila 2017. godine, nakon čega upisuje master akademske studije na Fakultetu tehničkih nauka, smer Inteligentni sistemi.