

Primjena metoda strojnog učenja u predviđanju kretanja vrijednosti burzovnog indeksa

Botunac, Ive

Master's thesis / Diplomski rad

2018

Degree Grantor / Ustanova koja je dodijelila akademski / stručni stupanj: **University of Zadar / Sveučilište u Zadru**

Permanent link / Trajna poveznica: <https://um.nsk.hr/um:nbn:hr:162:863920>

Rights / Prava: [In copyright](#) / [Zaštićeno autorskim pravom.](#)

Download date / Datum preuzimanja: **2024-10-19**



Sveučilište u Zadru
Universitas Studiorum
Jadertina | 1396 | 2002 |

Repository / Repozitorij:

[University of Zadar Institutional Repository](#)



zir.nsk.hr



DIGITALNI AKADEMSKI ARHIVI I REPOZITORIJ

Sveučilište u Zadru

Odjel za ekonomiju

Sveučilišni diplomski studij Menadžment (jednopedmetni – izvanredni)



Ive Botunac

**Primjena metoda strojnog učenja u predviđanju kretanja
vrijednosti burzovnog indeksa**

Diplomski rad

Zadar, 2018.

Sveučilište u Zadru

Odjel za ekonomiju

Sveučilišni diplomski studij Menadžment (jednopedmetni – izvanredni)

Primjena metoda strojnog učenja u predviđanju kretanja vrijednosti
burzovnog indeksa

Diplomski rad

Student:

Ive Botunac

Mentorica:

izv. prof. dr. sc. Anita Peša

Komentor:

doc. dr. sc. Ante Panjkota

Zadar, 2018.



Izjava o akademskoj čestitosti

Ja, **Ive Botunac**, ovime izjavljujem da je moj **diplomski** rad pod naslovom **Primjena metoda strojnog učenja u predviđanju kretanja vrijednosti burzovnog indeksa** rezultat mojega vlastitog rada, da se temelji na mojim istraživanjima te da se oslanja na izvore i radove navedene u bilješkama i popisu literature. Ni jedan dio mojega rada nije napisan na nedopušten način, odnosno nije prepisan iz necitiranih radova i ne krši bilo čija autorska prava.

Izjavljujem da ni jedan dio ovoga rada nije iskorišten u kojem drugom radu pri bilo kojoj drugoj visokoškolskoj, znanstvenoj, obrazovnoj ili inoj ustanovi.

Sadržaj mojega rada u potpunosti odgovara sadržaju obranjenoga i nakon obrane uređenoga rada.

Zadar, 10. srpnja 2018.

SAŽETAK

Predviđanja na tržištu kapitala jedan su od izazova kako za trgovce tako i za istraživače zbog svoje nelinearne i nestabilne strukture. S razvojem računalnih znanosti i informacijskih tehnologija nastaju novi pristupi temeljeni na strojnom učenju za poboljšanje procesa vezanih za trgovanje na tržištu kapitala. Ove nove metode prvenstveno su vezane za razvoj modela za predviđanje trendova pojedinih dionica ili burzovnih indeksa. U ovom radu korištene su povratne neuronske mreže (RNN) za predviđanje vrijednosti burzovnih indeksa (Dow Jones i NASDAQ) i za predviđanje kretanja trenda burzovnih indeksa u obliku binarne klasifikacije (porast/pad ili gore/dolje). U navedenim zadacima korišteno je pet različitih kombinacija ulaznih varijabli za model RNN: zaključna vrijednost indeksa, zaključna vrijednost indeksa s tehničkim indikatorima, zaključna vrijednost indeksa s naslovima Reuters portala novosti, zaključna vrijednost indeksa s FinSentS sentiment podacima i zaključna vrijednost indeksa s tehničkim indikatorima, naslovima Reuters portala novosti i s FinSentS sentiment podacima. Svi provedeni eksperimenti dali su rezultate usporedive s rezultatima dostupnih, relevantnih znanstvenih članaka. Najbolji rezultat u predviđanju vrijednosti burzovnog indeksa i predviđanja kretanja trenda u obliku gore/dolje dobiveni su kombinacijom RNN sa zaključnom vrijednosti burzovnog indeksa i FinSentS sentiment podacima. Ovi rezultati otvaraju prostor za brojne druge primjene koje koriste tehnike strojnog učenja za razvoj strategija trgovanja i sustava za automatizirano trgovanje.

Ključne riječi

predviđanja na tržištima kapitala, povratne neuronske mreže, tehnička analiza, sentiment-analiza

Applying machine learning methods to predict the movement of stock market index

ABSTRACT

Stock market predictions constitute one of the challenges for traders and explorers due to its nonlinear and unstable structure. With the developments in computer science and information technology, new approaches are arising based on machine learning to improve the processes related to stock market trading. These novel methods are primarily directed towards the development of models for the prediction of trends of the particular stock values or stock indexes. In this paper, recurrent neural networks (RNN) have been used for stock index values forecasting (Dow Jones and NASDAQ) and for predicting stock indexes trends in the form of binary classification (rising/falling or up/down). Five different combinations of input variables for RNN were used in mentioned tasks: closing index values, closing index values with technical indicators, closing index values with Reuters News titles classes, closing index values with FinSentS sentiments, and closing index values with Reuters News titles classes and FinSentS sentiments. All conducted experiments produced results comparable with the results in the accessible, relevant scientific articles. Best results in forecasting stock index values and predicting future trends in the form of up/down classes were obtained by combining RNN with closing values and FinSentS sentiments. These results are opening a space for numerous other applications which are using machine learning techniques in the tasks for developing of the trading strategy and systems for automated trading.

Keywords

stock markets forecast, recurrent neural network, technical analysis, sentiment analysis

SADRŽAJ

1. UVOD.....	1
2. TEORIJSKI OKVIR RADA	3
2.1. Strojno učenje.....	3
2.1.2. Područja strojnog učenja	4
2.1.3. Zadatci strojnog učenja.....	5
Regresija (engl. Regression).....	5
Klasifikacija (engl. Classification)	6
Klasteriranje (engl. Clustering)	7
Asocijacijska pravila (engl. Association rules)	7
Smanjenje dimenzionalnosti (engl. Dimensionality reduction)	8
Predviđanje (engl. Forecasting).....	8
2.1.4. Metode strojnog učenja	8
Stroj s potpornim vektorima (engl. Support Vector Machine).....	8
Stabla odlučivanja (engl. Decision Trees).....	9
Algoritam k-srednjih vrijednosti (engl. K-Means Clustering)	9
Naivan Bayesov klasifikator (engl. Naive Bayes Classifier)	10
Metode ansambla (engl. Ensemble methods).....	10
2.2. Umjetne neuronske mreže	10
2.2.1. Biološki živčani sustav	11
2.2.2. Model neurona.....	12
2.2.3. Aktivacijske funkcije.....	13
2.2.4. Proces učenja mreže	15
Funkcije pogreške (engl. Cost function).....	16

Algoritam najstrmijeg spusta (engl. Gradient Descent)	17
Algoritam povratne propagacije (engl. Backpropagation)	17
2.3. Financijski vremenski nizovi.....	19
2.3.1. Tržište kapitala	19
2.3.2. Koncept financijskih vremenskih serija	19
2.3.3. Tehnička analiza	21
2.3.4. Fundamentalna analiza	21
2.3.5. Sentiment-analiza	22
2.3.6. Hipoteza efikasnog tržišta	22
2.3.7. Bihevioralne financije.....	23
3. METODOLOGIJA, MODELI I OPIS PODATAKA	24
3.1. Metodologija.....	24
3.2. Modeli predviđanja.....	24
3.2.1. Povratne neuronske mreže.....	25
3.2.2. Čelija s dugoročnom memorijom	26
3.2.3. Specifikacija modela predviđanja.....	27
3.3. Prikupljanje i definiranje podataka.....	27
3.3.1. Burzovni indeksi.....	28
3.3.2. Tehnički indikatori	28
3.3.3. Reuters novosti	32
3.3.4. FinSentS	33
3.4. Obrada podataka	34
3.4.1. Ulazne i izlazne varijable	34
3.4.2. Skaliranje podataka.....	35

3.4.3. Podjela podataka.....	36
4. EKSPERIMENTI, REZULTATI I DISKUSIJA.....	38
4.1. Mjerni pokazatelji.....	38
4.1.1. Mjere kvalitete izvedbe za eksperimente regresijske analize.....	38
4.1.2. Mjere kvalitete izvedbe za eksperimente klasifikacije.....	39
4.1.3. Analiza statističkog značaja.....	40
4.2. Proces treniranja modela.....	41
4.2.1. Parametri modela predviđanja.....	41
4.2.2. Odabir parametra modela predviđanja.....	41
4.3. Rezultati eksperimenata.....	43
4.3.1. Rezultati regresijskog modela.....	43
4.3.2. Rezultati klasifikacijskog modela.....	47
5. ZAKLJUČAK.....	50
POPIS LITERATURE.....	51
POPIS SLIKA.....	57
POPIS TABLICA.....	58

1. UVOD

Tržište kapitala danas je važan pokazatelj gospodarskog rasta i razvoja zemlje. Tvrtke koje su sudionici tog tržišta od procesa trgovanja na burzama prikupljaju sredstva kako bi unaprijedile svoje procese, tehnologiju i infrastrukturu za daljnji razvoj. Upravo zbog važnosti tržišta kapitala postoji velik interes vezan za istraživanja budućih kretanja cijena kako bi se time ostvarili novčani prinosi.

Kretanja na tržištu kapitala su nelinearna i nestabilna što je uzrok složene prirode burza na kojima se odvija trgovanja, čime je prilično teško predvidjeti buduća kretanja. Postoje brojni čimbenici koji uzrokuju kolebanja kretanja cijena kao što su političke situacije, djelovanja poduzeća, određene gospodarske aktivnosti i drugi neočekivani događaji koji mogu uslijediti. S ovim tvrdnjama vidi se da je trgovcima i investitorima vrlo teško donijeti odluke o kupnji i ulaganju stoga se oni koriste određenim analizama kako bi mogli donijeti odluku.

Razvojem računalnih znanosti i informacijskih tehnologija javljaju se nove metode koje se koriste strojnim učenjem kao jednim od predstavnika umjetne inteligencije ne bi li se time pomoglo investitorima i trgovcima u procesima vezanim za trgovanja na tržištu kapitala. Korištenjem metoda strojnog učenja, prema rezultatima nekih od istraživanja (Liu *et al.*, 2017; Weng, 2017), vidljivo je da se kretanja na tržištu kapitala mogu predviđati do određenog stupnja točnosti što je nerijetko i iznad 70 %.

Cilj ovog rada je istražiti postojanje mogućnosti predviđanja kretanja buduće vrijednosti burzovnog indeksa. Burzovni indeks se uzima kao jedan od glavnih statističkih pokazatelja koji omogućuje prikaz stanja tržišta kapitala i njegovih promjena. Za potrebe provođenja eksperimenta koriste se povratne neuronske mreže (engl. *Recurrent Neural Network*) koje su se pokazale izrazito korisnim kod predviđanja podataka koji dolaze u vremenskim serijama.

Rad je podijeljen na nekoliko poglavlja gdje se od uvoda postupno razrađuje tematika prema cilju istraživanja.

Prvo poglavlje obrađuje teorijski okvir gdje se pokrivaju osnove strojnog učenja kao i financijska analiza. Detaljnije su obrađene umjetne neuronske mreže koje se kao osnova koriste za modele istraživanja ovog rada.

U drugom poglavlju opisuje se metodologija istraživanja s osvrtom na korišteni model predviđanja. Opisuju se kombinacije korištenih podataka koje se koriste u provođenju eksperimenta istraživanja.

U trećem poglavlju prikazani su rezultati provedenih eksperimenata te su objašnjene korištene mjere za kvalitetu izvedbe u slučaju regresije i klasifikacije. Također, diskutira se o dobivenim rezultatima i uspoređuje ih se s dobivenim rezultatima drugih autora.

2. TEORIJSKI OKVIR RADA

Ovim poglavljem rada daje se opći uvod u strojno učenje kako bi se u kasnijim poglavljima mogla lakše razumjeti primjena korištenih modela predviđanja. Posebna pozornost je posvećena detaljnijem opisu umjetnih neuronskih mreža (engl. *Artificial Neural Networks*), njihovih karakteristika i procesa učenja. Drugi dio teorijskog okvira rada posvećen je pojmu financijskih vremenskih nizova što je i predmet istraživanja. Tim djelom daje se uvod u analize tih nizova kako bi se lakše razumjela primjena modela predviđanja u rješavanju zadatka predviđanja financijskih vremenskih nizova.

2.1. Strojno učenje

Strojno učenje (engl. *Machine Learning*) je grana umjetne inteligencije (engl. *Artificial Intelligence*) koja u samoj osnovi izučava kreiranje algoritama sa sposobnošću učenja temeljem podataka ili kroz interakciju s okolinom. Rezultat tih napora su porodice algoritama koje rješavaju probleme iz različitih domena ljudskih djelatnosti kao što su financije, medicina, pravni sektor, industrijska proizvodnja i maloprodaja (The Royal Society, 2017). Za strojno učenje dodatno se može reći da je to programiranje računala ka optimizaciji kriterija uspješnosti pomoću primjera podataka ili prethodnog iskustva. Postoji model određen nekim parametrima, a učenje je u tome kontekstu izvršavanje računalnog programa kako bi se ti parametri optimizirali pomoću spomenutih primjera ili prethodnog iskustva (Alpaydin, 2009).

Marsland (2009) navodi da je tek proteklim desetljećem strojno učenje prepoznato u multidisciplinarnom pristupu s ostalim znanostima navodeći pri tome biologiju, statistiku, matematiku i fiziku, ali da se strojno učenje najviše promatra u području umjetne inteligencije čime se svrstava u polje računalnih znanosti.

Svoju primjenu strojno učenje danas nalazi u rješavanju brojnih problema gdje se to dokazalo učinkovitijim i bržim usporedo s ljudskim mogućnostima rješavanja istog problema. U objavljenom izvješću koje je izdao The Royal Society (2017) navode se neka od značajnih područja primjene strojnog učenja:

Sustavi preporuke (engl. *Recommender systems*) su u osnovi sustavi koji na temelju prethodnih izbora preporučuju određene usluge ili proizvode. Sustavi za preporuke koriste obrasce potrošnje i izražene preferencije potrošača kako bi predvidjeli koji će proizvod ili

uslugu preporučiti potrošaču. Takvi se sustavi rašireno koriste u mrežnim trgovinama koje su orijentirane ka maloprodaji.

Organiziranje informacije (engl. *Organising information*) vezano je za sustave koji pomažu u davanju rezultata prilikom upita upisanih u internetske tražilice. Ovdje spadaju i sustavi za otkrivanje neželjenih podataka što se opisuje primjerom detekcije neželjene elektroničke pošte (engl. *SPAM*).

Prepoznavanje govora (engl. *Voice recognition*): sustavi za obradu prirodnog jezika i prepoznavanje govora mogu se podudarati s uzorcima zvukova proizvedenih u ljudskom govoru riječima ili izrazima koje su već susreli. Nakon identificiranja upotrijebljenih riječi, ovi sustavi ih mogu prevesti u tekst ili izvršiti naredbe. Ovakvi sustavi danas sve bolje prepoznaju govor omogućivši da brojni pametni telefoni i drugi uređaji dolaze opremljeni virtualnim osobnim asistentima.

Računalni vid (engl. *Computer vision*) nalazi se integriran u sustavima koji mogu detektirati i analizirati slike te povezati numeričke ili simboličke informacije s tim slikama. U programima društvenih medija, prepoznavanje slika može se upotrebljavati za označavanje objekata ili osoba na fotografijama koje su prenesene na internetsku lokaciju.

Strojno prevođenje (engl. *Machine translation*) opisuje se kao računalne sustave koji mogu automatski pretvoriti tekst ili govor s jednog jezika na drugi. Danas se strojno prevođenje koristi u posebnim mobilnim aplikacijama, društvenim mrežama te u međunarodnim organizacijama koje trebaju reproducirati dokumente na velikom broju jezika.

Prepoznavanje uzoraka (engl. *Pattern recognition*) može se upotrijebiti za identifikaciju uzoraka u podacima koje ljudi ne bi mogli prepoznati. Zajednička primjena prepoznavanja uzoraka je u sustavima za otkrivanje prijevара povezanih s korištenjem kreditne kartice ili drugih platnih sustava. Primjerice, ako se za korisnika prikaže neobičan obrazac potrošnje, sustav može podići upozorenje.

2.1.2. Područja strojnog učenja

Strojno učenje obično se može podijeliti na dva glavna područja: nadzirano učenje (engl. *Supervised Learning*) i nenadzirano učenje (engl. *Unsupervised Learning*) koje je opisano u daljnjem tekstu poglavlja. Iako je ovo najzastupljenija podjela, kod nekih autora dodatno se

pronalaze podjele područja strojnog učenja na pojačavajuće učenje (engl. *Reinforcement Learning*) (Marsland, 2009; P. Murphy, 2012), aktivno učenje (engl. *Active learning*), učenje prijenosom (engl. *Transfer Learning*) (Yang, Hanneke i Carbonell, 2013), definiranje značajki (engl. *Features engineering*) (Nargesian *et al.*, 2017) te računalna teorija učenja (engl. *Computational learning theory*) (Blum, 2007). Kako nadzirano i nenadzirano strojno učenje spadaju pod najzastupljeniju podjelu strojnog učenja, te se nadzirano strojno učenje koristi u provođenju testiranja ovog rada, dodatno se opisuju te dvije podjele.

Nadzirano učenje je proces učenja gdje je cilj na osnovi ulaza x dobiti izlaz y trenirajući model sa zadanim ulazno-izlaznim parovima $M = (x_1, y_1, \dots, x_n, y_n)$. Ovdje se koristi M kao oznaku trening seta podataka, a n kao oznaku broja trening primjera. (P. Murphy, 2012).

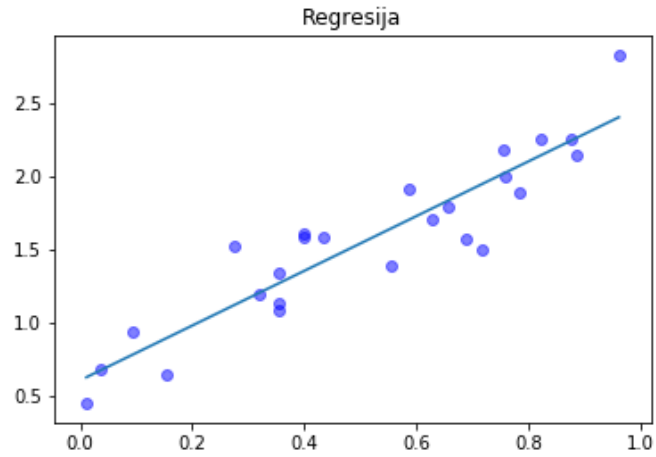
Kod nenadziranog učenja upotrebljava se drugačiji pristup jer u trening setu koji se koristi za treniranje modela nisu poznati izlazi y , već samo ulazne vrijednosti x , time je $M = (x_1, \dots, x_n)$. Stoga u ovoj vrsti strojnog učenja modeli pokušavaju identificirati sličnost između ulaza i na taj način ih grupirati (Bonnin, 2017).

2.1.3. Zadatci strojnog učenja

Danas postoje brojni zadatci koji se rješavaju primjenom strojnog učenja, a kao glavne može se istaknuti regresiju i klasifikaciju što je ujedno sadržano u predmetu istraživanja ovoga rada. Kako bi se steklo bolje razumijevanje i šira slika o strojnom učenju, u nastavku se uz regresiju i klasifikaciju navode i ostali važni zadatci strojnog učenja.

Regresija (engl. Regression)

Regresija se obično koristi kada se na temelju ulazne vrijednosti x želi predviđati kontinuiranu vrijednost izlaza y . Na Slici 1. može se vidjeti primjer gdje se postavlja regresijska funkcija na odnosu između dvije varijable koje u ovom slučaju predstavljaju ulaz x i izlaz y .

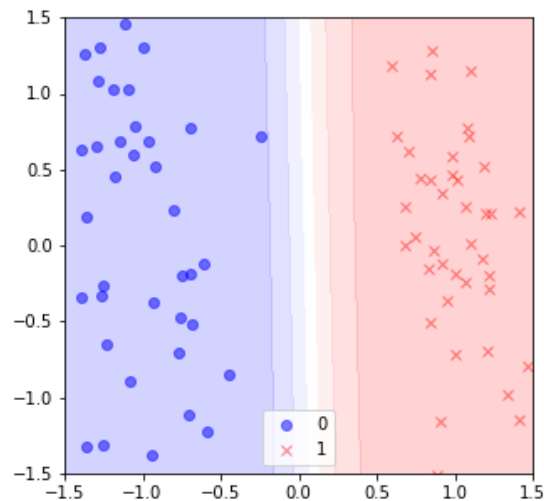


Slika 1. Linearna regresija

Izvor: izradio autor (2018)

Klasifikacija (engl. Classification)

Za razliku od regresije, kod klasifikacije se ne predviđa kontinuirano vrijednost, već se ulaz x svrstava u jednu od klasa y . Ako se za primjer uzme binarna klasifikacija, tada vrijednost y može biti samo 0 ili 1 što je prikazano na Slici 2.

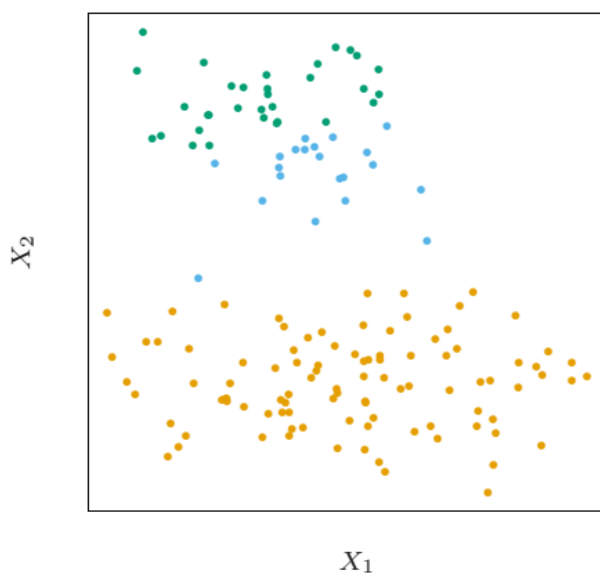


Slika 2. Binarna klasifikacija

Izvor: izradio autor (2018) prema Hastie, Tibshirani i Friedman (2001:13)

Klasteriranje (engl. Clustering)

Klasteriranje se odnosi na grupiranje ili segmentiranje objekata u podskupove ili klasterne, tako da su oni unutar svakog klastera blisko povezani jedan s drugim za razliku od objekata dodijeljenih drugim različitim skupinama. Objekt se može opisati skupom mjerenjima ili njegovim odnosom s drugim objektima. Osim toga, cilj ponekad može biti organizirati klasterne u prirodnu hijerarhiju (Hastie, Tibshirani i Friedman, 2001).



Slika 3. Simulirani podatci u ravnini, grupirani u tri klase

Izvor: Hastie, Tibshirani i Friedman (2001:502)

Na Slici 3. prikazuje se simulacija podataka koji su grupirani u tri skupine putem jednog od algoritama koji se koristi za rješavanje ovog zadatka strojnog učenja.

Asocijacijska pravila (engl. Association rules)

Cilj ovog zadatka strojnog učenja je pronaći zajedničke vrijednosti varijabli $X = (X_1, X_2, \dots, X_n)$ koje se najčešće pojavljuju u promatranoj bazi podataka. Time se često primjenjuju binarne vrijednosti $X \in \{0, 1\}$ što se u ovome kontekstu naziva „analiza tržišne košarice“ (Hastie, Tibshirani i Friedman, 2001). Asocijacijsko pravilo je implikacija oblika $X \rightarrow Y$ gdje je X prethodnica, a Y je posljedica pravila (Alpaydin, 2009).

Smanjenje dimenzionalnosti (engl. Dimensionality reduction)

Smanjivanje dimenzionalnosti koristi se kod rješavanja kompleksnih problema, primjerice klasifikacije ili regresije koji se zasnivaju na velikom broju ulaznih atributa. Korištenjem smanjenja dimenzionalnosti može se znatno ubrzati izvršavanje modela, poboljšati rezultate modela i olakšati rad s podacima (Marsland, 2009).

Prema Alpaydin (2010), javljaju se dvije metode koje se koriste za smanjivanje dimenzionalnosti, a to su: odabir značajki (engl. *Feature selection*) i ekstrakcija značajki (engl. *Feature extraction*).

Predviđanje (engl. Forecasting)

Predviđanje je specifičan zadatak vezan za strojno učenje, a najizraženiji je u kontekstu predviđanja vremenskih nizova. Kod predviđanja vremenske serije postoji skup podataka koji pokazuju kako se nešto razlikuje tijekom vremena i stoga se želi predvidjeti kako će se podatci u budućnosti razlikovati odnosno kretati. Ovo je prilično težak zadatak strojnog učenja, ali koristan je u bilo kojem području istraživanja gdje se podatci pojavljuju u vremenskim nizovima (Marsland, 2009).

2.1.4. Metode strojnog učenja

Kako bi se pristupilo rješavanju prethodno navedenih zadataka strojnog učenja, mora se odabrati određena metoda koja će biti najbolja za rješavanje tog zadatka. U istraživanju ovog rada koristit će se umjetne neuronske mreže koje se posebno detaljno opisuju u točki 2.2. Osim umjetnih neuronskih mreža, u nastavku se navode neke od danas najpoznatijih metoda.

Stroj s potpornim vektorima (engl. Support Vector Machine)

Stroj s potpornim vektorima je metoda nadziranog strojnog učenja koja se koristi većinom za rješavanje klasifikacijskih zadataka. Ova metoda radi na način da neki skup označenih podataka dijeli u određene klase koristeći pri tome liniju razdvajanja (Tandel, 2017).

Stroj s potpornim vektorima pokušava maksimizirati udaljenost između različitih klasa koje se nalaze u skupu podataka. Ovu udaljenost naziva se marginom iz čega slijedi što je margina

veća, to je manja pogreška generalizacije klasifikatora. Važna karakteristika ove metode strojnog učenja je da se rješenje temelji samo na onim podatkovnim točkama koje se nalaze na rubovima margina nazvanih potpornim vektorima (Tandel, 2017).

Stabla odlučivanja (engl. Decision Trees)

Stabla odlučivanja također spadaju pod metodu nadziranog strojnog učenja te se uglavnom primjenjuju za rješavanje klasifikacijskih zadataka. Ovaj model klasificira podatke u skupu podataka prolazeći kroz strukturu upita iz korijena dok ne dođe do lista, što predstavlja jednu klasu. Korijen predstavlja atribut koji igra glavnu ulogu u klasifikaciji, a list predstavlja klasu. Jednostavno rečeno, stablo odlučivanja je stablo u kojem svaki čvor predstavlja izbor između brojnih alternativa, a svaki list predstavlja odluku.

Mohammed, Khan i Mohammed Bashier (2016) navode sljedeće korake u razvrstavanju podataka:

1. Stavljaju sve trening primjere u korijen stabla.
2. Dijele trening primjere na temelju odabranih atributa.
3. Odabiru attribute pomoću nekih statističkih mjera.
4. Rekurzivno dijeljenje nastavlja sve dok ne ostane nijedan trening primjer ili dok ne ostane nijedan atribut.

Algoritam k-srednjih vrijednosti (engl. K-Means Clustering)

Algoritam k-srednjih vrijednosti spada pod vrstu nenadziranog strojnog učenja čiji je cilj grupirati podatke u grupe na osnovi zadanog broja grupa od strane korisnika. Ova metoda strojnog učenja koristi iterativno usklađivanje kako bi se dobio konačan rezultat u vidu podjele zadanog seta podataka na grupe minimiziranjem udaljenosti od k-srednjaka (centroida) (Mohammed, Khan i Mohammed Bashier, 2006).

Naivan Bayesov klasifikator (engl. Naive Bayes Classifier)

Ova metoda strojnog učenja izrazito je popularna za rješavanje klasifikacijskih zadataka. Naivan Bayesov klasifikator temelji se na Bayesovu teoremu te on predviđa kolika je vjerojatnost da zadani podatak pripada određenoj klasi.

Klasa koja ima najveću vjerojatnost smatra se konačno najvjerojatnijom klasom kojoj podatak pripada. Naivni se naziva iz razloga što polazi od pretpostavke da je vrijednost jednog atributa u nekoj klasi neovisna od vrijednosti ostalih atributa (Shalev-Shwartz i Ben-David, 2014).

Metode ansambla (engl. Ensemble methods)

Metode ansambla spadaju pod tehniku strojnog učenja koja kombinira nekoliko osnovnih metoda strojnog učenja kako bi se time stvorio optimalan model. Ovakvo kombiniranje metoda strojnog učenja je izrazito učinkovito jer se iskorištavanjem raznolikosti korištenih metoda može smanjiti ukupnu grešku čime se povećava točnost modela (Zhang i Ma, 2012). Postoji nekoliko različitih pristupa u korištenju metode ansambla od čega se izdvajaju:

„Bagging“ pristup koji koristi metodu ponavljanog uzrokovanja iz skupa podataka kako bi se time generirali slučajni skupovi podataka.

„Boosting“ pristup koji se temelji na ideji da u slučaju više korištenih modela strojnog učenja jedan pokriva ono područje u kojem neki drugi ne daje tako dobre rezultate. Ovaj pristup stavlja veću težinu na pogrešno klasificirane rezultate, a istodobno smanjuje težinu ispravno klasificiranih.

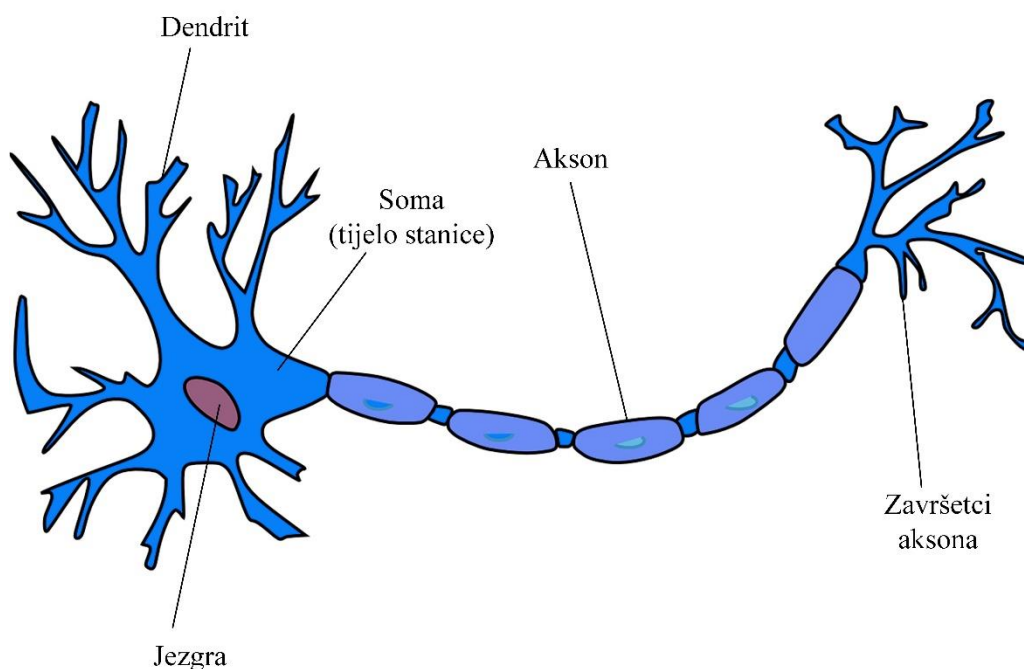
2.2. Umjetne neuronske mreže

Ovom točkom izdvajaju se umjetne neuronske mreže kao posebna cjelina iz razloga što su korištene u provođenju eksperimenta ovog rada i time je potrebno detaljnije opisati njih kao metodu strojnog učenja. Umjetne neuronske mreže spadaju pod vrstu nadziranog strojnog učenja te je njihovo nastajanje inspirirano biološkim živčanim sustavom (McCulloch i Pitts, 1943). Da bi se stekla intuicija o umjetnim neuronskim mrežama, prvo se ukratko opisuje sam biološki živčani sustav te potom model neurona koji je izveden iz neurona opisanog u

biološkom živčanom sustavu. Nakon ovog uvoda upisuje se sam proces treniranja umjetnih neuronskih mreža i aktivacijske funkcije koje se koriste.

2.2.1. Biološki živčani sustav

Kako bi se bolje razumjelo način na koji funkcioniraju umjetne neuronske mreže, mora se prvo upoznati živčani sustav koji je poslužio kao inspiracija u razvoju ove metode strojnog učenja. Ljudski mozak sastoji se od gotovo 10 milijardi međusobno povezanih živčanih stanica koje se nazivaju neuroni (Delbelo Bašić, Čupić i Šnajder, 2008). Za kreiranje modela neurona dovoljno je promatrati stvarni neuron kroz sljedeće dijelove: soma, što je samo tijelo stanice, dendrite, koji su mala vlakna oko some te aksone, koji su duže vlakno. Na Slici 4. prikazana je ilustracija biološkoga neurona.



Slika 4. Biološki neuron

Izvor: izradio autor (2018) prema (Kriesel, 2005:17)

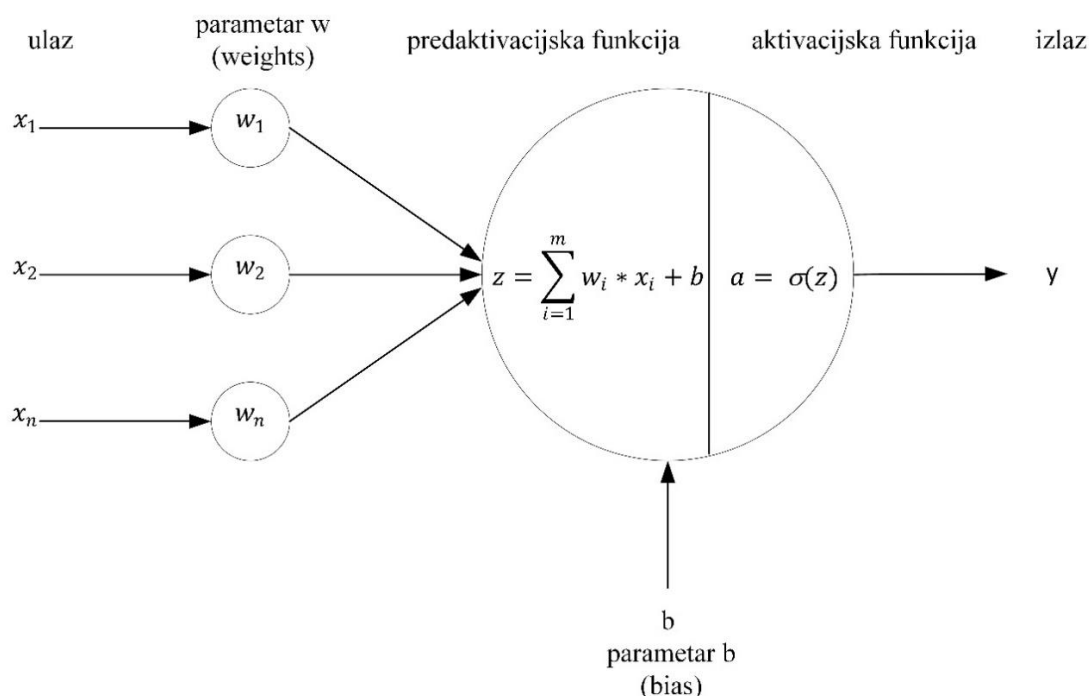
Dolazni signali iz drugih neurona prenose se posebnim vezama nazvanim sinapse koje predstavljaju male pukotine između dendrita susjednih neurona. Ako je jedan neuron pobuđen, što dovodi do putovanja akcijskog potencijala aksonom, kada taj potencijal dođe do završetaka, on uzrokuje otpuštanje različitih neurotransmitera u prostor sinapsi, ti neurotransmiteri se primaju receptorima na drugom kraju i uzrokuju postsinaptički potencijal. Svi postsinaptički

potencijali se zbrajaju i čine novi potencijal koji se prenosi aksonom, čime se postupak nastavlja. Djelovanje neurotransmitera može biti takvo da nastoje pobuditi neuron ili ga deaktivirati.

2.2.2. Model neurona

Svaki čvor ili neuron unutar neuronske mreže ima određeni broj ulaznih kanala i jedan izlazni kanal. Za svaki od neurona može se reći da obavlja dvije operacije, nad ulazima obavlja se sumiranje što se dodatno naziva predaktivacijom koja predstavlja ulaz u tzv. prijenosnu funkciju što se naziva aktivacijom. Prema Haykin (2009) ovime se može identificirati tri osnovna elementa modela neurona, te je vidljiv na Slici 5.:

- [1] Skup sinapsi ili veza od kojih je svaka okarakterizirana težinom koja čini njezinu vlastitu snagu. Signal x_i ulazne sinapse koja je spojena na neuron množi se s težinom (engl. *Weight*) w_i
- [2] Zbrajanje svih ulaznih signala x_i pomnoženih s odgovarajućim težinama w_i
- [3] Funkcija aktivacije za ograničavanje amplitude izlaznog signala neurona u kojoj se smanjuje dopušteni raspon amplitude izlaznog signala do neke konačne vrijednosti.



Slika 5. Model neurona

Izvor: izradio autor (2018) prema Haykin, (2009:11)

Model neurona također primjenjuje dodatni parametar b (*engl. Bias*) koji predstavlja pomak te ima efekt na povećavanje ili smanjivanje neto ulaza aktivacijske funkcije. Matematički model neurona prikazan na Slici 5. može se prikazati jednadžbama:

za predaktivaciju

$$z = \sum_{i=1}^m w_i * x_i + b \quad (1)$$

te za aktivaciju slijedi

$$a = \sigma(z) \quad (2)$$

gdje w_i predstavlja težinski parametar ulaznog i -tog signala, time je x_i pripadna ulazna vrijednost. Oznakom m označava se broj ulaznih parova \vec{x} i y koji se koriste kao trening podatci. Nakon izračunane predaktivacijske vrijednosti z koristi se aktivacijske funkcije σ kako bi se dobilo konačan izlaz a .

2.2.3. Aktivacijske funkcije

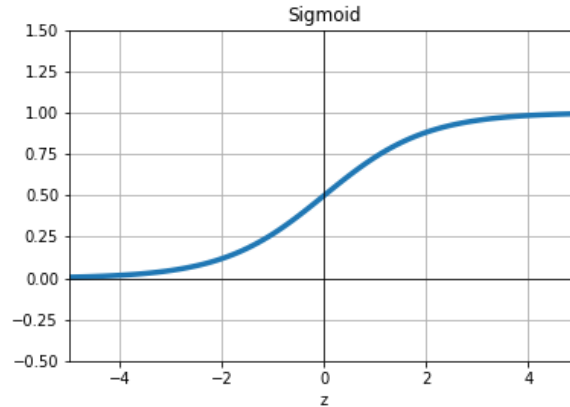
Aktivacijske funkcije su iznimno važne kod umjetnih neuronskih mreža. One u osnovi odlučuju hoće li se određeni neuron u mreži aktivirati i prenijeti informaciju ka drugim neuronima u mreži (Ramachandran, Zoph i Le, 2017). Najjednostavnija aktivacijska funkcija je step funkcija (*engl. step function*) koja vraća samo dvije vrijednosti: nulu ili jedinicu na temelju toga je li ulazni podatak manji ili veći od određene granice, pa se ponekad naziva i funkcija praga (*engl. Threshold function*). Među najpopularnijim aktivacijskim funkcijama danas se nalaze sigmoidne, tanh i relu funkcije.

Sigmoidna funkcija je također poznata i kao logistička aktivacijska funkcija. Funkcija uzima neki realan broj z te ga zatim skalira unutar granica 0 i 1. Većinom se koristi u izlaznom sloju neuronske mreže gdje je krajnji cilj predvidjeti vjerojatnost temeljem binarne klasifikacije na nulu i jedinicu (Goodfellow, Bengio i Courville, 2016). Sigmoidna funkcija pretvara velike negativne brojeve na 0 i velike pozitivne brojeve na 1 čime se u konačnici dobiva da rezultat klasifikacije može biti jedna od dvije moguće vrijednosti.

Matematički se može prikazati jednadžbom i grafički na Slici 5.:

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (3)$$

Izgled amplitude sigmoidne funkcije prikazan je na Slici 6.



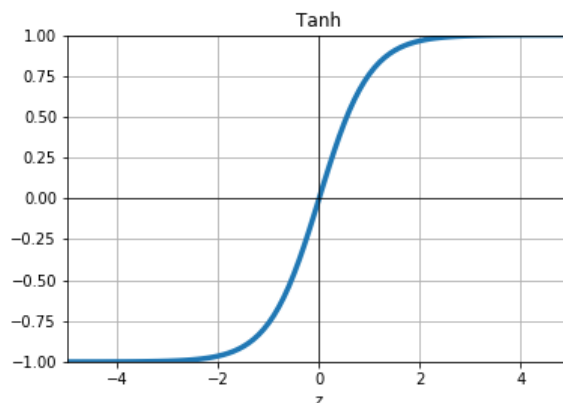
Slika 6. Sigmoid funkcija

Izvor: izradio autor (2018)

Tanhova funkcija je aktivacijska funkcija slična sigmoidnoj funkciji koja vraća vrijednosti između -1 i 1. Za razliku od sigmoidne funkcije, ova funkcija može mapirati u pozitivna i negativna područja. Time su veliki negativni brojevi skalirani prema -1, a veliki pozitivni brojevi skalirani su prema 1. Matematički se može prikazati kao:

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (4)$$

Izgled amplitude tanh funkcije prikazan je na Slici 7.



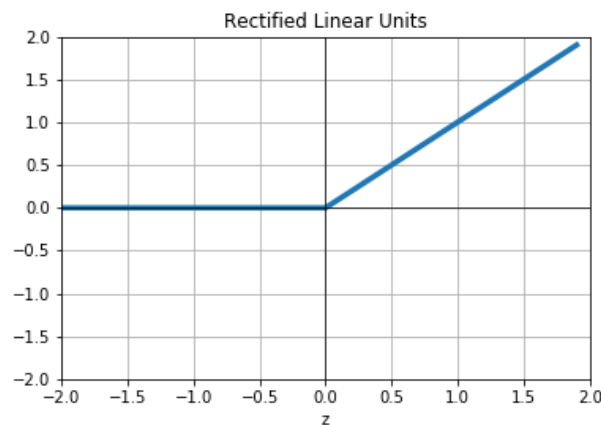
Slika 7. Tanh funkcija

Izvor: izradio autor (2018)

ReLU funkcija jedna je od najčešće korištenih aktivacijskih funkcija unutar neuronskih mreža. Glavna prednost korištenja ReLU funkcije u odnosu na druge aktivacijske funkcije je u tome što ne aktivira sve neurone u isto vrijeme. Ako je ulaz negativan, pretvara ga u nulu i neuron se ne aktivira (Maas, Hannun i Ng, 2013). To znači da se samo određeni broj neurona aktivira čime mreža postaje učinkovita i jednostavnija za izračunavanje. ReLU aktivacijska funkcija prikazana je sljedećom matematičkom formulom:

$$\sigma(z) = \max(0, z) \quad (5)$$

gdje funkcija za ulaz $z < 0$ vraća 0, a za $z > 0$ vraća z . Upravo je ovim načinom neuronska mreža brža i dolazi do znatno manjih zasićenja prilikom izračuna (Dan-Ching, 2017). Izgled amplitude ReLU funkcije prikazan je na Slici 8.

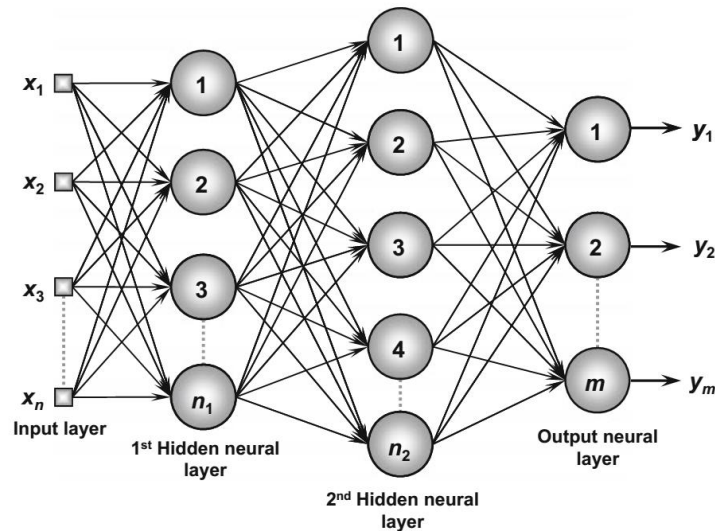


Slika 8. ReLU funkcija

Izvor: izradio autor (2018)

2.2.4. Proces učenja mreže

Proces učenja neuronske mreže definiran je mogućnošću da mreža sama sebi prilagođava vrijednosti parametra; pri tome se misli na težinske parametre w i parametar b . Parametri se kroz proces učenja postupno prilagođavaju sve dok je razlika izlaza iz mreže i stvarnog izlaza koji se nalazi u skupu podataka za treniranje minimalna (Kriesel, 2005).



Slika 9. Primjer višeslojne umjetne neuronske mreže

Izvor: Silva *et al.* (2017:23)

Na Slici 9. prikazana je jednostavna neuronska mreža koja se sastoji od ulaznog sloja, dva skrivena sloja i jednog izlaznog sloja. Ulazni sloj neuronske mreže predstavlja ulazne značajke na osnovi kojih se izračunava izlaznu značajku. Unutar skrivenih slojeva provode se predaktivacijske i aktivacijske funkcije opisane u točki 2.2.2. ovog rada. Završno izlazni sloj daje rezultat na osnovi svih provedenih izračuna.

Funkcije pogreške (engl. Cost function)

Učinak procesa učenja mreže može se mjeriti tzv. funkcijama pogreške. Dvije najpoznatije mjere koje se koriste su prosječna kvadratna greška (engl. *Mean Squared Error*) i unakrsna entropija (engl. *Cross Entropy*). U svome istraživanju Golik, Doetsch i Ney (2013) dolaze do zaključka da su bolji rezultati ostvareni primjenom unakrsne entropije kao funkcije pogreške.

Ove funkcije pogreške mjere razliku između očekivane izlazne vrijednosti i dobavne izlazne vrijednosti. Što je ta vrijednost bliža nuli, to je naš proces učenja postigao zadovoljavajuće vrijednosti treniranih parametara (Nielsen, 2015). Funkciju prosječne kvadratne greške matematički se opisuje sljedećim izrazom:

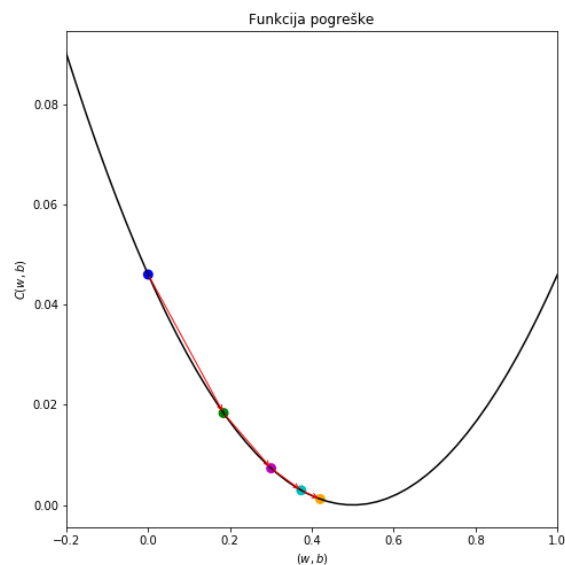
$$C(w, b) = \frac{1}{2m} \sum_{i=1}^m \|y(x^{(i)}) - a\|^2 \quad (6)$$

dok se unakrsnu entropiju prikazuje izrazom:

$$C(w, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} * \log a^{(i)} + (1 - y^{(i)}) * \log(1 - a^{(i)})] \quad (7)$$

Algoritam najstrmijeg spusta (engl. Gradient Descent)

Algoritam gradijentnog spusta jedan je od češće korištenih načina optimizacije parametra unutar neuronske mreže. On spada pod optimizacijske algoritme koji se koriste za minimiziranje neke funkcije koristeći iterativne pomake u smjeru pronalaženja globalnog minimuma odnosno najniže vrijednosti (Goodfellow, Bengio i Courville, 2016).



Slika 10. Gradijentni spust

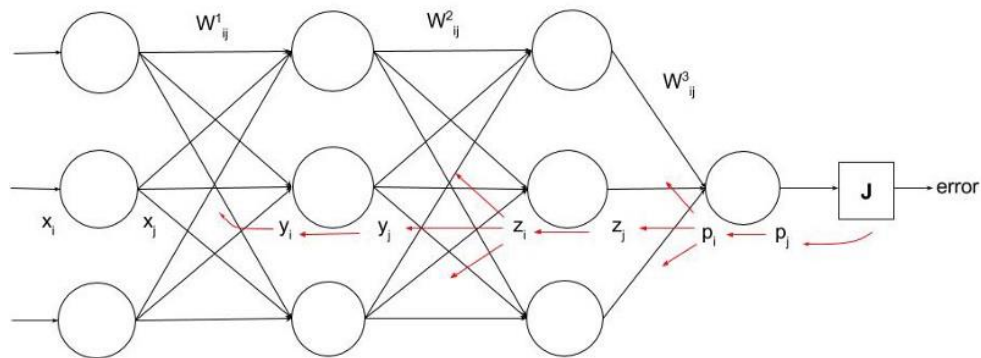
Izvor: izradio autor, (2018)

Na Slici 10. je prikazano minimiziranje konveksne funkcije pogreške koristeći gradijentni spust.

Algoritam povratne propagacije (engl. Backpropagation)

Algoritam povratne propagacije pruža način računanja gradijentnog pada funkcije pogreške unutar neuronske mreže. On ne spada pod metodu optimizacije već samo označava proces

primjene metode optimizacije kroz neuronsku mrežu kako bi se podešavali njezini parametri što je prikazano Slikom 11.



Slika 11. Algoritam povratne propagacije

Izvor: Kapur i Khazan (2016)

Uzimajući u obzir algoritam povratne propagacije, proces učenja može se svesti na dvije osnovne faze:

- [1] faza označava prosljeđivanje ulazne vrijednosti kroz mrežu uzimajući u obzir izračune koje provode modeli neurona sve do njezina kraja
- [2] faza započinje izračunom funkcije pogreške te propagacijom unatrag kroz mrežu kako bi se korištenjem optimizacijskog algoritma podešavale vrijednosti.

Završni rezultat treba uračunati u jednadžbu za prilagodbu novonastalih parametara w i b koristeći unaprijed definiranu vrijednost stope učenja. Iz toga slijedi podešavanje parametra w :

$$w_i^+ = w_i - \alpha * \frac{\partial C}{\partial w_i} \quad (8)$$

te zatim podešavanje parametra b :

$$b_i^+ = b_i - \alpha * \frac{\partial C}{\partial z} \quad (9)$$

2.3. Financijski vremenski nizovi

Financijski vremenski nizovi čine drugu komponentu predmeta istraživanja ovog rada pri tome navodeći strojno učenje kao prvu komponentu. U ovome dijelu cilj je pružiti osnovne informacije o tržištu kapitala te dati uvod u analize koje se provode nad tim tržištem posebno ističući tehničku, fundamentalnu i sentiment-analizu.

2.3.1. Tržište kapitala

Tržište kapitala definira se kao tržište na kojem se provodi proces ekonomske razmjene ili trgovine između kupca i prodavača, kao što je trgovanje devizama, obveznicama ili dionicama. Darskuviene (2010) navodi da tržište kapitala igra ključnu ulogu u gospodarstvu poticanjem gospodarskog rasta, utječući na ekonomsko djelovanje aktera, koji utječu na ekonomsko blagostanje. To se postiže financijskom infrastrukturom u kojoj subjekti s novcem dodjeljuju sredstva onima koji imaju potencijalno produktivnije načine ulaganja tih sredstava. Struktura financijskog tržišta omogućuje kupcima i prodavačima određivanje cijene i vrijednosti financijskih potraživanja ili pak željene stope povrata na različite vrste financijske imovine. Nadalje, financijsko tržište nudi likvidnost ulagačima kroz mogućnost pretvaranja financijske imovine u likvidna sredstva.

Važan dio tržišta kapitala čine burze kao mjesto gdje se nalaze prodavači i kupci kako bi ostvarili trgovanje. Sve aktivnosti koje se provode na burzama vidljive su svim sudionicima što značajno može utjecati na cijenu. Dionice se smatraju likvidnim ako se uoče značajne aktivnosti trgovanja. Što je više aktivnosti, lakše je pronaći kupca kada netko pokuša prodati i obrnuto.

2.3.2. Koncept financijskih vremenskih serija

Na tržištu kapitala, podatci financijskih vremenskih serija definirani su kao nizovi ponovljenih promatranih varijabli. Može se nabrojiti primjerice cijene dionica, tečajne liste, povrat obveznica i cijene robe mjereno u ujednačenim vremenskim razmacima. Osnovne varijable financijskih vremenskih serija uključuju vremensku oznaku, cijenu otvaranja, najvišu cijenu, najnižu cijenu, cijenu zatvaranja i volumen trgovanja. Cijene se prate u određenoj vremenskoj frekvenciji i time se stvaraju vremenske serije. Osnovna značajka financijskih vremenskih

serija je visoka učestalost pojedinačnih vrijednosti (Arlt i Arltová, 2001). Primjer jedne financijske vremenske serije burzovnog indeksa prikazan je na Slici 12.



Slika 12. Kretanje burzovnog indeksa „Dow Jones“

Izvor: izradio autor (2018)

Financijske vremenske serije služe kao osnova za provođenje analiza na tržištima kapitala. Prema Tsayu (2010), analiza financijskih vremenskih serija odnosi se na primjenu teorije i prakse vrednovanja imovine tijekom vremena. Navodi da je to empirijska disciplina, ali kao i ostala znanstvena polja, teorija čini temelje za donošenje zaključaka.

Nadalje, analize se provode kako bi se mogla ostvariti neka određena predviđanja u smislu budućih kretanja tih serija. Predviđanje se odnosi na vjerojatne čimbenike koji mogu utjecati na buduće poslovanje. Također identificira trendove kako bi se donijela odluka o budućnosti financijskih ulaganja. Međutim, predviđanje financijskih vremenskih serija je nelinearno i dinamičko zbog nestacionarnosti. Tradicionalno, dobiveni podatci vremenske serije ne sadrže dovoljno podataka za razumijevanje budućih trendova, stoga je vrlo teško predviđati buduće kretanje samo na osnovi financijskih vremenskih serija. Kao primjer može se navesti neuspjeh ekonomista u predviđanju gospodarske krize koja se dogodila 2008. godine (Krugman, 2009).

2.3.3. Tehnička analiza

Tehnička analiza je proučavanje kretanja tržišta, prvenstveno korištenjem grafikona u svrhu predviđanja budućih kretanja cijena. Prema Murphyju (1999), postoje tri glavna principa na kojima je bazirana tehnička analiza:

1. U cijeni su sadržani svi fundamentalni faktori. Tehnički analitičari smatraju da sve što bi eventualno moglo utjecati na cijenu pri tome navodeći fundamentalne faktore, političke, psihološke ili neke druge faktore zapravo je već ukomponirano u samoj cijeni. Dakle oni zaključuju da je proučavanje kretanja cijene sve što je analitičaru potrebno.
2. Cijene se kreću u trendovima. Koncept trenda je apsolutno neophodan za tehnički pristup. Cilj je na tržištu identificirati trendove u ranoj fazi njihova razvoja u svrhu trgovanja u smjeru tih trendova. U stvari, većina tehnika korištena u ovom pristupu slijedi trend, što znači da je njihova namjera identificirati i pratiti postojeće trendove. Postoji korelacija pretpostavci da se cijene kreću u trendovima, trend kretanja je – vjerojatnije je da će se nastaviti nego da se preokrenu.
3. Povijest se ponavlja. Većina elemenata tehničke analize i proučavanje tržišnih kretanja ima veze s proučavanjem ljudske psihologije. Primjerice obrasci grafikona koji su identificirani i kategorizirani u posljednjih stotinu godina sadržavaju određene obrasce koji se pojavljuju na cjenovnim tablicama. Budući da su ti obrasci dobro funkcionirali u prošlosti, pretpostavlja se da će i dalje dobro funkcionirati u budućnosti.

Tehnička analiza koristi modele i pravila trgovanja temeljena na promjenama cijena i volumena, kao što su indeks relativne čvrstoće (engl. *Relative Strength Index*), prosjeci kretanja (engl. *Moving Averages*), regresije, ciklusi ili klasično prepoznavanjem uzoraka grafikona. Tehnička analiza analizira cijenu, količinu i druge tržišne informacije, dok fundamentalna analiza razmatra stvarne činjenice tvrtke, tržišta, valute ili robe.

2.3.4. Fundamentalna analiza

Fundamentalna analiza smatra se jednim od najlakših načina vrednovanja poduzeća čiji je glavni cilj otkrivanje stvarne trenutne vrijednosti tvrtke. Za ciljeve fundamentalne analize Baresa, Bogdan i Ivanovic (2013) navode predviđanje buduće dobiti, dividendi i rizika kako bi se izračunala stvarna vrijednost dionica. Analiza se sastoji od proučavanja gospodarstva u

cjelini gdje se koriste određeni makroekonomski faktori kao i relevantne informacije o industriji u kojoj djeluje sama tvrtka. U obzir prilikom analiziranja ulazi cjelokupni učinak tvrtke, njezina financijska izvješća, uključujući i sve najnovije vijesti o tvrtki. Na temelju svega treba zaključiti je li tržište ispravno procijenilo sve informacije prilikom formiranja cijene dionica. Investitor treba uzeti u obzir sve dijelove financijskih izvještaja, uključujući dobit, imovinu, prihode i rashode, napraviti usporednu analizu po godini, napraviti usporednu analizu prema određenim industrijskim standardima, primijetiti određene trendove u njihovu ponašanju i na temelju svega toga pravilno procijeniti vrijednost dionice (Petrusheva i Jordanoski, 2016).

Prema (Thomsett, 1998), fundamentalna analiza može biti vrijedan alat ako se koristi za ostvarivanje dugoročnog profita, ali ne i za praćenje svakodnevnog kretanja cijena dionica, tržišne reakcije na vijesti ili glasine, te privremene popularnosti jedne industrijske grupe nad drugom.

2.3.5. Sentiment-analiza

Sentimentalna analiza označava proces računalnog utvrđivanja prenosi li tekst objavljen na portalima ili društvenim mrežama pozitivno, negativno ili neutralno mišljenje korisnika. Ova vrsta analize također se koristi za praćenje i analizu socijalnih fenomena, za uočavanje potencijalno opasnih situacija i određivanje općeg raspoloženja na internetu koje je uzrokovano reakcijom na razne vijesti (Pimprakar, Ramachandran i Senthilkumar, 2017).

Kao glavni zadatak Pawar, Jawale i Kyatanavar (2016) spominju istraživanje problema klasifikacije subjektivnosti i klasifikacije sentimenta kod provođenja klasifikacije nad nekim tekstom. Klasifikacija subjektivnosti je polje koje identificira sadržava li dan tekstualni dokument činjenične podatke ili informacije o sebi. Tada je razvrstavanje po sentimentu odgovorno za kategoriziranje mišljenja u pozitivno mišljenje ili negativno.

2.3.6. Hipoteza efikasnog tržišta

Ovo je teorija koja navodi da su sve relevantne informacije uključene u formiranju cijene dionice ne bi li time investitorima osigurali ostvarivanje prosječnog prinosa na tržištu (Barbić, 2010). Glavna posljedica toga bila bi da je gotovo nemoguće nadmudriti cjelokupno tržište. Teorija se temelji na pretpostavci da su sve relevantne informacije javno i lako dostupne svim

investitorima te da investitori djeluju racionalno (Fama, 1970). Prirodna posljedica ove teorije je da su financijske vremenske serije uvijek nepredvidljive. Teorija slučajnog hoda (engl. *Random Walk Theory*) navodi da su sve financijske vremenske serije statistički ekvivalentne nizu sasvim slučajnih koraka, a pokušaj da se predviđanja temelje na tehničkoj analizi rezultirat će neuspjehom. Ako je to istina, onda bi općenito bilo nemoguće da algoritamski sustavi trgovanja i predviđanja budućih kretanja budu profitabilni.

Logika teorije slučajnog hoda prema Malkielu (2003) je ta da ako je protok informacija neometan i informacije se odmah odražavaju u cijenama dionica, onda će sutrašnja promjena cijene odražavati samo sutrašnje vijesti i bit će neovisna o promjenama cijena danas. Ali vijest je po definiciji nepredvidljiva, pa stoga promjene cijena moraju biti nepredvidljive i slučajne. Kao rezultat, cijene u potpunosti odražavaju sve poznate informacije, pa čak i neinformirani investitori koji kupuju raznoliki portfelj na tablici cijena koje daje tržište dobit će jednaku stopu povrata kao i onu koju bi postigli stručnjaci.

2.3.7. Bihevioralne financije

Bihevioralne financije istražuju psihološki aspekt u donošenju financijskih odluka te su sasvim nov pristup području financija. Brajković i Peša (2015) navode da „Tradicionalno se smatra da pri odlučivanju u uvjetima neizvjesnosti ljudi oblikuju stavove u skladu sa zakonima vjerojatnosti i da se pri tome vode isključivo maksimizacijom osobnih interesa. Međutim, ne treba zanemariti i značaj intuicije, prema kojem zaključivanje i donošenje odluka dolazi spontano, bez iscrpnog razmišljanja i napora.“

Kako je ovo područje relativno novije, samim time pruža se zanimljiva alternativa klasičnim financijama. Klasične financije pretpostavljaju da su tržišta kapitala učinkovita, investitori su racionalni i dugoročno nije moguće nadmašiti tržište. Psihološka načela ponašanja uključuju, između ostalog, predrasude, prekomjerno povjerenje, emocije i određene društvene sile. Time bihevioralne financije lako objašnjavaju zašto je pojedinac poduzeo određenu odluku, ali nije lako pronaći objašnjenje o tome kako će taj isti pojedinac donositi buduće odluke. Prema hipotezi efektivnog tržišta, s obzirom na to da svatko ima pristup istim informacijama, gotovo je nemoguće pobijediti tržište, jer su cijene dionica zapravo učinkovite, odražavajući sve što investitori znaju. Nasuprot tome, bihevioralne financije pretpostavljaju da su u nekim okolnostima financijska tržišta informativno neučinkovita (Biräu, 2012).

3. METODOLOGIJA, MODELI I OPIS PODATAKA

Ovim poglavljem opisuju se različiti pristupi i metode koje su korištene u svrhu provođenja eksperimenta i dobivanja rezultata. Detaljno se opisuju korišteni podatci i načini na koje su oni generirani ili odakle su preuzeti.

3.1. Metodologija

Glavni zadatak u provođenju istraživanja ovog rada je testirati utjecaj različitih kombinacija i vrsta ulaznih podataka kako bi se time moglo predviđati buduća kretanja na tržištu kapitala primjenom modela povratnih neuronskih mreža. Koristi se kombinacija sva tri navedena pristupa u prethodnom poglavlju pri analiziranju financijskih vremenskih nizova, a to su redom: tehnička, fundamentalna i sentiment-analiza. Testiranja su provedena nad dva američka burzovna indeksa: NASDAQ i Dow Jones Industrial Average koja generalno opisuju ukupna kretanja na američkom tržištu kapitala.

Iz ovoga se definiraju dvije osnovne empirijske metode pri provođenju testiranja. Prvom metodom na temelju različitih kombinacija ulaznih podataka predviđa se sutrašnja vrijednost burzovnog indeksa što je vezano za regresijsko predviđanje budućih kretanja. Drugom metodom također na temelju kombinacije različitih ulaznih podataka predviđa se sutrašnje kretanje trenda. Tu se postavlja problem binarne klasifikacije u vidu ako je sutrašnji trend u rastu, dodjeljuje mu se binarna vrijednost 1, a ako je trend u padu, dodjeljuje mu se binarna vrijednost 0.

U nastavku poglavlja obrađuju se korišteni modeli za predviđanje te se detaljno opisuju podatci koji su korišteni u spomenutim kombinacijama.

3.2. Modeli predviđanja

Za model predviđanja u svrhu provođenja istraživanja odabrana je jedna vrsta neuronskih mreža koja se naziva povratne neuronske mreže (engl. *Recurrent Neural Network – RNN*) s ćelijama s dugoročnom memorijom (engl. *Long-Short Term Memory – LSTM*). Ova vrsta neuronskih mreža koristi se za obradu sekvencijalnih podataka gdje spadaju i financijski vremenski nizovi

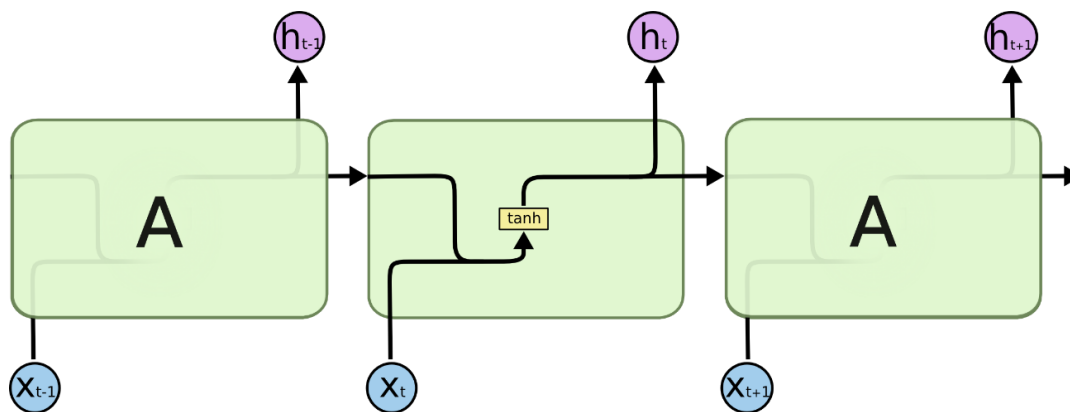
(Widegren, 2017). Njihova upotreba u svrhu predviđanja kretanja na tržištu kapitala može se pronaći u brojnim radovima (Hansson, 2017; Nelson, Pereira i De Oliveira, 2017) te u mrežnim repozitorijima (Olah, 2015; Munoz, 2018).

Odabir ovog modela racionalan je izbor jer je istraživanje ovog rada usredotočeno na predviđanja kretanja financijskih vremenskih serija gdje su se povratne neuronske mreže pokazale znatno boljim nego standardne unaprijedne neuronske mreže (Nelson, Pereira i De Oliveira, 2017; Widegren, 2017). Općenito neuronske mreže bilo da se radi o povratnim ili unaprijednim kod predviđanja financijskih vremenskih serija ostvaruju bolje rezultate od drugih metoda strojnog učenja kao što su stroj s otpornim vektorima ili stabla odluke (Patel *et al.*, 2015b, 2015a). Ipak, potrebno je napomenuti da su najbolji rezultati u dostupnoj literaturi ostvareni u istraživanju koje je proveo Weng (2017) koristeći se metodom ansambla koja kombinira više različitih metoda strojnog učenja.

3.2.1. Povratne neuronske mreže

Povratne neuronske mreže predstavljaju često korištenu metodu za obradu sekvencijalnih podataka poput strojnog prevođenja, predviđanja unosa teksta, prepoznavanja govora i predviđanja kretanja financijskih vremenskih serija (Widegren, 2017). Kod standardnih unaprijednih neuronskih mreža (engl. *Feedforward Neural Network*) podatci se kreću u samo jednom smjeru od ulaznog sloja prema izlaznom dok se kod povratnih može reći da podatci kruže kroz petlju (Donges, 2018). Povratna neuronska mreža unutar svakog čvora koristi određenu internu memoriju (engl. *Context Unit*) koja uzima informacije s prethodnih izlaza i standardno prosljeđenu ulaznu informaciju. Može se reći da povratna neuronska mreža svoj izračun generira na temelju sadašnje informacije i nedavne prošlosti. Upravo zbog ovoga svojstva, povratne neuronske mreže pokazale su se iznimno korisne kod sekvencijalnih podataka koji se mijenjaju s vremenom jer će u obzir uzeti prethodno naučene informacije iz danih vremenskih serija (Nielsen, 2017).

Na Slici 13. je prikaz povratne neuronske mreže gdje se vidi na koji način se informacije prosljeđuju iz jednog čvora u drugi.

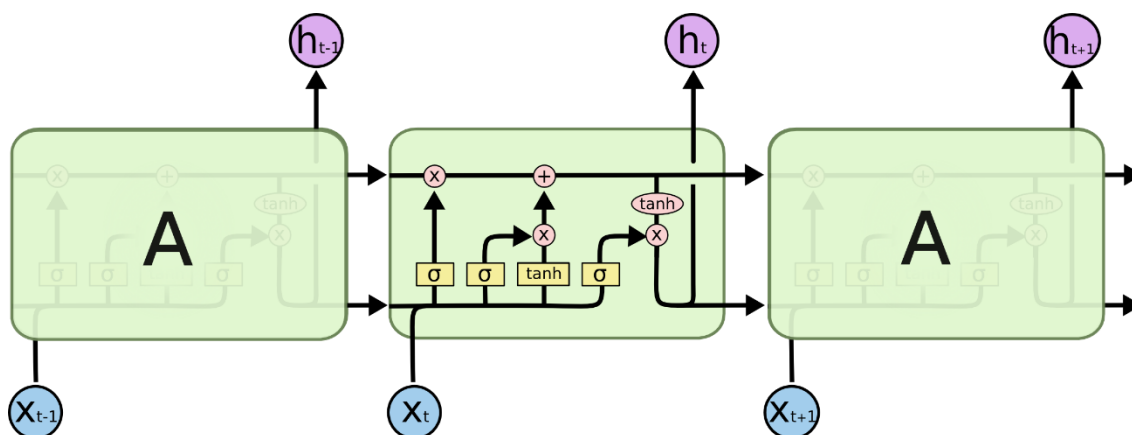


Slika 13. Povratna neuronska mreža

Izvor: Olah (2015), URL: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

3.2.2. Čelija s dugoročnom memorijom

Kod povratnih neuronskih mreža javljali su se određeni problemi prilikom izračuna te je izrazito teško uspješno provesti proces njezina treniranja (Nielsen, 2017). U rješavanju ovog problema razvijena je posebna ćelija s dugoročnom memorijom koju su izvorno predstavili Hochreiter i Schmidhuber (1997). Čelije s dugoročnom memorijom imaju sposobnost uklanjanja ili dodavanja informacija pojedinoj ćeliji što je reguliraju posebne unutarnje strukture nazvane vrata (engl. *Gates*) (Olah, 2015). Kao rezultat toga mogu se blokirati određeni signali koji su u običnim povratnim neuronskim mrežama uzrokovali probleme.



Slika 14. Čelija s dugoročnom memorijom

Izvor: Olah (2015), URL: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Na Slici 14. je prikazan primjer ćelije s dugoročnom memorijom gdje se vide ulazi i ilustracijski prikazi operacija koji se provode nad informacijama. Unutar ćelije je sadržano troje vrata koja služe kao zaštita i kontrola promjena stanja ćelije.

3.2.3. Specifikacija modela predviđanja

Za izradu modela predviđanja korišten je programski jezik Python (Python, 2018) s razvojnim okruženjem Keras (Keras, 2018). Keras je razvojno okruženje visoke razine apstrakcije namijenjeno za lako i brzo prototipiranje modela umjetnih neuronskih mreža. Dodatno za rukovanje podacima koristi se programska biblioteka Pandas (Pandas, 2018) s kojom su se podatci pripremali u za to prikladne tablice kako bi se potom isti mogli koristiti u modelima predviđanja. Model predviđanja izvorno su razvili i oblikovali autori u svrhu provođenja testiranja ovog rada gdje su za inspiraciju poslužili neki od izvora (Brownlee, 2016b, 2016a; Munoz, 2018).

3.3. Prikupljanje i definiranje podataka

Nakon definiranja modela predviđanja za rješavanje navedenih zadataka koji su predmet istraživanja ovog rada potrebno je prikupiti podatke. U samom istraživanju koriste se različiti skupovi podataka kao što su zaključne vrijednosti burzovnih indeksa, odabrani tehnički indikatori, naslovi s Reuters portala koji se dodatno obrađuju i komercijalni sentiment podatci s mrežne baze podataka FinSents.

Podatci su prikupljeni s nekoliko različitih izvora i to iz vremenskog perioda od 1. 1. 2013. godine do 7. 5. 2018. godine što ukupno čini skup od 1330 skupova podataka. S portala Alpha Vantage (2018) prikupljeni su podatci vezani za zaključne vrijednosti burzovnih indeksa, s novinskog portala Reuters (2018) prikupljeni su naslovi novinskih članaka koji su se koristili za fundamentalnu analizu dok su s portala FinSentS (2018) prikupljeni podatci vezani za elemente sentiment-analize.

3.3.1. Burzovni indeksi

Burzovni indeksi služe za opisivanje performansi tržišta dionica ili određenog dijela tog tržišta. Sastoje se od hipotetskog portfelja vrijednosnih papira te se izračunavaju kao aritmetičke sredine, jednostavne ili ponderirane. Formiraju se na način da kompanija koja ima znatno veći utjecaj na određeno tržište isto tako ima i veću težinu u samom indeksu (Šlibar, 2009). Primjeri indeksa dionica su NASDAQ, S&P 500 i Dow Jones Industrial Average.

Dow Jones Industrial Average (DJIA) je najcitiraniji indeks dionica na svijetu. Promjene u indeksu često se smatraju reprezentativnim za cijelu burzu (Shoven, 2000). Indeks je otkrio Charles H. Dow 26. svibnja 1896. i tad se indeks sastojao od 12 dionica što se do današnjeg dana povećalo na ukupno 30 dionica najvećih američkih tvrtki. Za ovaj indeks se kaže da je cjenovno ponderiran indeks što znači da dionice s višim cijenama imaju veću težinu odnosno utjecaj na sam indeks.

NASDAQ Composite Indeks (IXIC) je pokrenut 1971. godine te danas obuhvaća više od 3.000 vrijednosnih papira što uključuje dionice, investicijske fondove, fondove za nekretnine i brojne druge vrste (Krein i W. Smith, 2014). Indeks se izračunava temeljem metodologije ponderirane tržišne kapitalizacije čime je vrijednost indeksa jednaka ukupnoj vrijednosti pondera udjela svakog od konstitutivnih vrijednosnih papira, pomnožen s posljednjom cijenom svake pojedine vrijednosnice (Investopedia, n.d.).

3.3.2. Tehnički indikatori

Tehnički indikatori ili pokazatelji su matematički izračuni temeljeni na cijeni, količini ili otvorenom interesu nekog vrijednosnog papira, primjerice dionice ili burzovnog indeksa. Njih koriste uglavnom aktivni trgovci kako bi mogli analizirati kretanja cijena i na temelju toga predviđati buduća kretanja (Investopedia, n.d.).

Brojni autori u svojim istraživanjima (Tsai, C.-F. i Wang, 2009; Patel *et al.*, 2015a; Weng, 2017) koriste razne kombinacije tehničkih indikatora ne bi li temeljem njih poboljšali performanse svojih modela. U daljnjem tekstu opisuju se korišteni tehnički indikatori kako bi se mogao vidjeti njihov utjecaj na konačne rezultate predviđanja. Dodatno se koristi oznaka „TI“ kako bi se skupno označilo sve navedene tehničke indikatore koji se koriste.

Jednostavni pomični prosjek (engl. *Simple Moving Average – SMA*) je jedan od najsvestranijih i najčešće korištenih svih tehničkih indikatora. Ovaj indikator koristi prosjeke zaključnih cijena određenog broja dana kako bi se kreirala krivulja i time olakšalo uočavanje trenda (Murphy, 1999).

Jednadžbom se prikazuje izračun jednostavnog pomičnog prosjeka za 10 dana:

$$SMA = \frac{P_n + P_{n-1} + \dots + P_{n-9}}{10} \quad (10)$$

čime se vidi da se zbrajaju sve zaključne cijene prethodnih 10 dana i ukupna vrijednost dijeli se s brojem 10.

Ponderirani pomični prosjek (engl. *Weighted Moving Average – WMA*) stavlja veći naglasak na nedavne promjene cijena nego jednostavni pomični prosjek. Cijena svakog dana ima svoju težinu ovisno o tome kako se nedavno dogodila što znači da se veća težina dodjeljuje novijim podacima (Murphy, 1999).

Izračun ponderiranog pomičnog prosjeka prikazuje se jednadžbom:

$$WMA = \frac{(10)P_n + (9)P_{n-1} + \dots + P_{n-9}}{n + (n - 1) + \dots + 1} \quad (11)$$

Moment (engl. *Momentum – MOM*) je indikator koji služi za mjerenje brzine promjene cijene i na taj način indirektno prati kretanje krivulje zaključne cijene (Colby, 2003).

Njegov izračun može se prikazati formulom:

$$M = P - P_n \quad (12)$$

gdje se izračunava kao razlika između zadnje zaključne cijene i zaključne cijene prije n dana.

Prosječno kretanje cijena konvergencija/divergencija (engl. *Moving Average Convergence Divergence – MACD*) se izračunava kao razlika između dva eksponencijalna pomična prosjeka gdje se u praksi najčešće uzima vrijednosti od 12 i 26 dana. Eksponencijalni pomični prosjeci

naglašavaju nedavne promjene cijena slično kao i jednostavni pomični prosjeci samo što se kod eksponencijalnog pomičnog prosjeka veća težina dodjeljuje novijoj vrijednosti (Person, 2004).

Ovaj tehnički indikator može se prikazati jednadžbama:

$$MACD = EMA(12) - EMA(26) \quad (13)$$

dok se eksponencijalni pomični prosjek prikazuje:

$$EMA = (P - E_p) * K + E_p \quad (14)$$

gdje se oznakom P označuje cijenu zatvaranja, E_p se koristi kao oznaku eksponencijalnog pomičnog prosjeka prethodnog perioda te oznaka K predstavlja konstantu jednaku $2/(n + 1)$ iz čega je n broj promatranih perioda (Person, 2004).

Stohastički oscilator %K i %D (engl. *Stochastic oscillator - %K %D*) su indikatori koji prate i uspoređuju cijenu zatvaranja s cjenovnim rasponom tijekom određenog razdoblja. Trenutačna cijena tada se izražava kao postotak tog raspona s 0 % što ukazuje na dno raspona i 100 % što ukazuje na gornju granicu raspona tijekom promatranog vremenskog razdoblja (Person, 2004).

Stohastički oscilator %K se izračunava putem jednadžbe:

$$\%K = 100 \frac{P - Pl_n}{Ph_n - Pl_n} \quad (15)$$

Gdje P označava cijenu zatvaranja trenutnog dana, Pl_n i Ph_n označavaju najnižu i najvišu cijenu tijekom promatranog perioda.

Potom slijedi jednadžba za stohastički oscilator %D:

$$\%D = EMA(\%K, 3) \quad (16)$$

Iz ovih izračuna dobiveni su rezultati za dvije krivulje od čega %K krivulja predstavlja samu vrijednost indikatora, a %D krivulja se izračunava kao eksponencijalni pomični presjek 3 vrijednosti %K krivulje.

Indeks relativne čvrstoće (*Relative Strength Index* – *RSI*) prikazuje trenutnu i povijesnu snagu ili slabost promatranog vrijednosnog papira temeljem zaključnih cijena nedavno promatranog perioda. Dodatno se ovaj indikator klasificira kao momentni oscilator koji mjeri brzinu i veličinu usmjerenih kretanja cijena. RSI je najčešće korišten u 14-dnevnom vremenskom okviru mjerenom na ljestvici od 0 do 100, s visokom i niskom razinom označenom na 70 i 30 (Romeu, 2001).

$$RSI = 100 - \frac{100}{1 + RS} \quad (17)$$

gdje je

$$RS = \frac{EMA(U, n)}{EMA(D, n)} \quad (18)$$

Za svako promatrano razdoblje izračunava se navedena promjena u povećanju prosjeka povećanja cijene U ili prosjeka smanjenja cijene s D što se prikazuje jednadžbama:

$$U = P_n - P_{n-1} \quad (19)$$

$$D = 0$$

$$D = P_{n-1} - P_n \quad (20)$$

$$U = 0$$

Williams %R oscilator (WILLR) prikazuje se kao razlika trenutne cijene te najviše i najniže cijene u promatranom periodu. Prilikom izračuna najčešće se odabiru razdoblja od 14 i 28 dana (Murphy, 1999).

Williams %R oscilator prikazuje se jednadžbom:

$$\%R = 100 \frac{P_{h(n)} - P}{P_{h(n)} - P_{l(n)}} \quad (21)$$

gdje oznake indeksa $h(n)$ i $l(n)$ prikazuju najvišu odnosno najnižu zaključnu cijenu promatranog razdoblja.

Indeks robnog kanala (*engl. Commodity Channel Index CCI*) služi za identifikaciju cikličkih kretanja na robnim tržištima. Često se koristi za otkrivanje odstupanja od kretanja cijena, time primjerice ako se vrijednost diže iznad brojke +100 signalizira kupovni signal, u suprotnom ako pada ispod brojke -100 signalizira prodajni signal (Colby, 2003).

Indeks se izračunava uzimanjem razlike između trenutne cijene i njezina jednostavnog pomičnog prosjeka, podijeljene s prosječnim apsolutnim odstupanjem cijene što se prikazuje:

$$CCI = \frac{P_t - SMA(P_t)}{(0,015)\sigma(P_t)} \quad (22)$$

gdje je

$$P_t = \frac{P_h + P_l + P_c}{3} \quad (23)$$

Iz čega je σ prosječno apsolutno odstupanje od P_t , te P_h označava najvišu cijenu, P_l najnižu i P_c označuje cijenu zatvaranja.

3.3.3. Reuters novosti

Za izvor novosti vezanih za američko tržište kapitala koristi se portal pod nazivom Reuters (Reuters, 2018). Njega se smatra relevantnim izborom, prema radovima drugih autora koji koriste naveden izvor novosti (Schumaker i Chen, 2009; Atkins, Niranjana i Gerding, 2018) te prema istraživanju portala Investopedia (2018).

S navedenog portala za svaki dan u prethodno navedenom odabranom vremenskom intervalu arhivirani su naslovi svih objavljenih novosti. Za dohvaćanje novosti razvijen je poseban program koristeći programski jezik Python (Python, 2018) koji vrši iteraciju nad portalom Reuters za svaki dan odabranog vremenskog intervala te time naslove svih novosti sprema u tablice.

Prilikom dohvaćanja naslova novosti koristeći se metodama za obradu prirodnog jezika (engl. *Natural Language Processing*) mjeri se pozitivnost, negativnost, neutralnost, subjektivnost i polarnost što se naziva sentiment-analizom.

Kako bi se provela ova mjerenja, razvijen je program koji koristi model VADER (engl. *Valence Aware Dictionary and sEntiment Reasoner*) koji su razvili Hutto i Gilbert (2014). Ovaj model zasnovan je na kombiniranju različitih rječnika koji za određene riječi engleskoga jezika imaju predefinirana mjerenja za svaku pojedinačnu riječ.

Leksički pristup promatra mjere svake riječi u rečenici te time pronalazi prosjek mjera za svaku rečenicu. Prednost leksičkih pristupa leži u činjenici da se ne treba trenirati model pomoću označenih podataka, budući da je prisutno sve što je potrebno za procjenu sentiment-rečenica (Calderon, 2017).

3.3.4. FinSentS

FinSentS je mrežna baza podataka koja prikuplja podatke od nekoliko tisuća internetskih stranica, blogova i poslovnih vijesti svakih 5 minuta u potrazi za informacijama gdje se spominje tražena tvrtka ili burzovni indeks. Nad svakim od prikupljenih članaka koji su prema određenim algoritmima FineSentS mrežne baze podataka ocijenjeni kao relevantni za predmet upita primjenjuje se sentiment-analiza kako bi se dobilo određene mjerne pokazatelje. Skupni sentiment svih članaka računa se kao prosječna vrijednost svih prikupljenih sentimenta (InfoTrie, n.d.).

Ovdje se radi o komercijalnoj bazi podataka gdje se za pristup mora pretplatiti kako bi se podatci mogli preuzeti. U svrhu provođenja testiranja izvršena je pretplata na navedenu mrežnu bazu podataka te su za prethodno spomenut odabrani vremenski interval dohvaćeni sljedeći podatci:

Sentiment-vrijednost (engl. *Sentiment Score*) je mjera kojom se označava na skali od 1 do 10 koliko je sentiment odabranog burzovnog indeksa na određeni dan. Vrijednosti od 1 do 3 označuju negativni sentiment, vrijednosti od 4 do 7 označuju neutralan sentiment dok vrijednosti od 8 do 10 označuju pozitivan sentiment.

Količina vijesti (engl. *News Volume*) je mjera kojom se označava količina vijesti objavljena i analizirana na određeni dan. Ova mjera je dobar pokazatelj preokreta trenda.

News Buzz pokazuje kolika je promjena standardne devijacije od periodičnih količina vijesti. Mjeri se u skali od 1 do 10.

Tablica 1. Prikaz vrijednosti preuzetih s mrežne baze podataka FinSentS za burzovni indeks Dow Jones

Datum	Sentiment-vrijednost	Količina vijesti	News Buzz
2013-01-01	7	26	6,8
2013-01-02	6,167	106	8,5
2013-01-03	6,45	86	8,1
2013-01-04	5,7	62	7,4
2013-01-05	5,25	13	4

Potrebno je dodatno napomenuti da se radi o komercijalnoj mrežnoj bazi podataka koja svoje algoritme ne otkriva javnosti. U Tablici 1. nalazi se prikaz nekoliko uzoraka preuzetih vrijednosti s FinSentS mrežne baze podataka.

3.4. Obrada podataka

Obradom podataka podrazumijeva se primjenu odgovarajućih tehnika u cilju transformacije podataka u oblik prihvatljiv za primjenu unutar odabranih modela strojnog učenja. U nastavku slijedi opis načina kojim su se izvorni podatci obradili i pripremili sukladno prethodno definiranom modelu povratnih neuronskih mreža s ćelijom s dugoročnom memorijom.

3.4.1. Ulazne i izlazne varijable

Prije samog početka provođenja daljnjih koraka u vidu obrade podataka moraju se definirati ulazne varijable na temelju kojih se predviđaju izlazne varijable. Kako bi se testiralo utjecaj različitih financijskih analiza definira se nekoliko kombinacija ulaznih varijabli kako bi se one mogle primijeniti u modelima predviđanja u svrhu provođenja eksperimenta.

Kako je ranije navedeno, pristupilo se rješavaju dva specifična zadatka u predviđanju vrijednosti finansijskih vremenskih serija. Prvi zadatak je predviđanje vrijednosti narednog dana pa se temeljem toga kao izlazna varijabla uzima upravo cijena narednog dana. Za ulazne varijable koriste se kombinacije prikazane u Tablici 2.

Tablica 2. Popis kombinacija ulaznih atributa za provođenje eksperimenata predviđanja vrijednosti burzovnih indeksa ili klase rasta/pada njihove vrijednosti

Kombinacija	Opis
K1	zaključna vrijednost
K2	zaključna vrijednost, TI
K3	zaključna vrijednosti, pozitivnost, negativnost, neutralnost, subjektivnost i polarnost (Reuters News)
K4	zaključna vrijednost, sentiment, volume, news buzz (FinSentS)
K5	zaključna vrijednost, TI, pozitivnost, negativnost, neutralnost, subjektivnost i polarnost (Reuters News), sentiment, volume, news buzz (FinSentS)

Kod rješavanja klasifikacijskog zadatka predviđanja kretanja sutrašnjeg trenda za izlaznu varijablu mora se generirati binarnu klasifikaciju trenda na temelju zaključnih vrijednosti. Za generiranje binarne klasifikacije trenda koristi se jednostavno pravilo: ako je *vrijednost* (i) veća od *vrijednosti* ($i-1$), onda je trend „1“ (u rastu), u suprotnom trend je „0“ (u padu). Isto kao i kod regresijskog zadatka, za ulazne varijable koristi se kombinacije podataka prikazanih u Tablici 2.

3.4.2. Skaliranje podataka

Vrlo važan korak u obradi podataka je skaliranje ulaznih vrijednosti odnosno ulaznih varijabli. Podatci koji su skalirani omogućuju da se modeli strojnog učenja znatno brže izvršavaju (Rechenthin, 2014). Cilj ove tehnike obrade podataka je skalirati ulazne vrijednosti unutar zadanih intervala kako bi se izbjegle velike vrijednosti ili velika odstupanja između različitih vrijednosti ulaznih varijabli. Ovom tehnikom, primjerice, može se sve vrijednosti skalirati u intervalu od -1 i 1 što se prikazuje sljedećom jednadžbom:

$$x_{ts} = \frac{x_t - x_{min}}{x_{max} - x_{min}} \quad (24)$$

U svrhu provođenja testiranja ovog rada zbog razlika u rasponu vrijednosti sve varijable u svim kombinacijama skalirale su se na vrijednosti u rasponu od -1 i 1. Primjer izvornih i skaliranih vrijednosti prikazan je u Tablici 3.

Tablica 3. Prikaz izvornih i skaliranih vrijednosti zaključne vrijednosti i vrijednosti preuzetih s mrežne baze podataka FinSentS za burzovni indeks Dow Jones

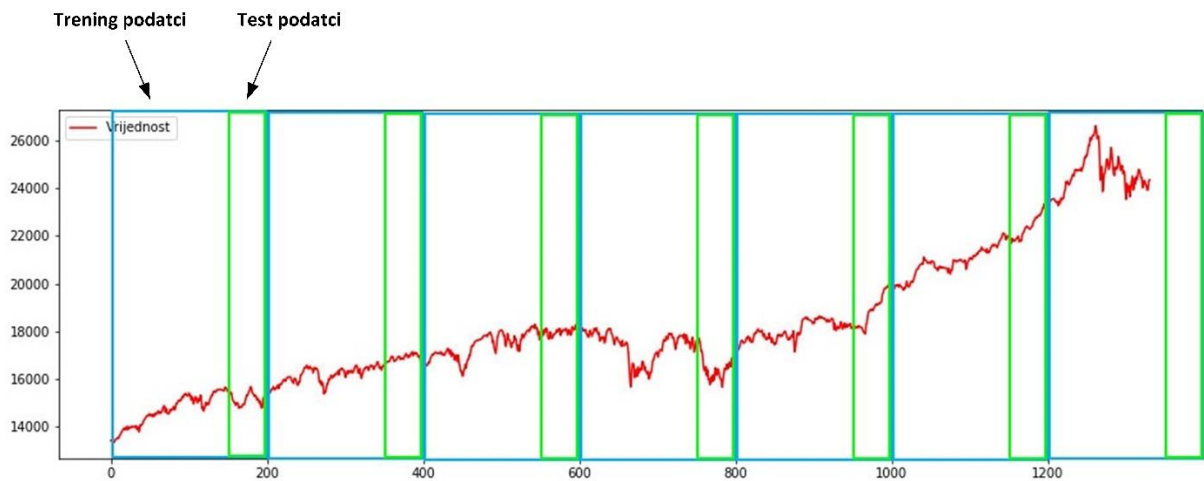
Datum	Zaključna vrijednost		Sentiment-vrijednost		Količina vijesti		News Buzz	
	Izvorno	Skalirano	Izvorno	Skalirano	Izvorno	Skalirano	Izvorno	Skalirano
2013-01-02	13412.5498	0.006299	6.167	0.6167	106	0.158209	8.500	0.8500
2013-01-03	13391.3604	0.004704	6.450	0.6450	86	0.128358	8.100	0.8100
2013-01-04	13435.2100	0.008004	5.700	0.5700	62	0.092537	7.400	0.7400
2013-01-07	13384.2900	0.004172	6.000	0.6000	128	0.191045	7.643	0.7643
2013-01-08	13328.8496	0.000000	6.583	0.6583	107	0.159701	6.833	0.6833

3.4.3. Podjela podataka

Pod preliminarnim radnjama prije provođenja eksperimenata također se podrazumijeva podjela podataka na skup podataka koji se koristi za treniranje modela i na skup koji se potom koristi za testiranje modela. Tradicionalno se koriste podjele u omjeru od 80 % za trening skup te 20 % za testni skup što može varirati ovisno o veličini ukupno dostupnih podataka.

Kod financijskih vremenskih serija javljaju se određeni problemi kod podjele podataka koji nastaju zbog specifičnosti samog predmeta istraživanja. Standardnom podjelom podataka može se nazvati onu podjelu gdje 80 % trening podataka čini vremenski niz od početka odabranog razdoblja što se kontinuirano nastavlja na 20 % testnog. Rechenthin (2014) navodi da prilikom standardne podjele može nastati problem prevelike zastupljenosti jedne klase ili obujma

vrijednosti unutar trening seta što rezultira ukupno lošijim rezultatima. To je osobito vidljivo kod kontinuirane serije podataka, pri čemu se temeljna struktura podataka može vremenom mijenjati zbog promjene tržišne dinamike.



Slika 15. Prikaz pomičnog prozora

Izvor: izradio autor (2018)

Kako bi se izbjegli navedeni problemi, koristi se tehnika podjele podataka koja se naziva klizni prozor (engl. *Sliding Window*) ako se radi o predviđanju buduće vrijednosti (Rechenthin, 2014; Weng, 2017). Ako se radi o klasifikacijskom problemu kod predviđanja budućeg trenda koristi se slična tehnika koja se naziva stratificirano uzorkovanje (engl. *Stratified Sampling*) (Rechenthin, 2014; Weng, 2017).

Pomoću ovih tehnika brine se o ravnomjernoj podjeli serije podataka na način da su trening i testni skup uzeti iz svih dijelova dostupnog skupa podataka što je prikazano na Slici 15.

4. EKSPERIMENTI, REZULTATI I DISKUSIJA

Definiranom metodologijom i modelima predviđanja dolazi se do provođenja eksperimenta kako bi se moglo vidjeti utjecaj različitih kombinacija ulaznih podataka na predviđanje buduće vrijednosti i kretanje trenda. Prije samih rezultata definiraju se mjerni pokazatelji koji se koriste sukladno rješavanju regresijskog, odnosno klasifikacijskog zadatka.

4.1. Mjerni pokazatelji

Kako bi se rezultati provedenog eksperimenta mogli interpretirati i razumjeti, potrebno je obraditi odabrane mjerne pokazatelje. U ovom radu obrađivani su problemi regresije i klasifikacije stoga se za svaki pojedini problem koriste različite mjere kvalitete izvedbe koje se posebno obrađuju.

4.1.1. Mjere kvalitete izvedbe za eksperimente regresijske analize

Za mjerenje uspješnosti regresijskih modela koriste se brojni mjerni pokazatelji koji definiraju odstupanje predviđene vrijednosti od stvarne vrijednosti. U radu se koriste srednja apsolutna pogreška (engl. *Mean Absolute Error – MAE*) i korijen srednje kvadratne pogreške (engl. *Root Mean Squared Error – RMSE*).

Srednja apsolutna pogreška je prosjek razlike između stvarnih vrijednosti i predviđenih vrijednosti. Ona daje mjeru koliko daleko su predviđene vrijednosti udaljene od stvarnih vrijednosti, ali bez smjera udaljenosti (Mishra, 2018). Srednju apsolutnu pogrešku izračunava se pomoću jednadžbe:

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (25)$$

Srednja kvadratna pogreška vrlo je slična srednjoj apsolutnoj pogreški. Razlika je u tome što srednja kvadratna greška uzima prosjek kvadrata razlike između stvarnih vrijednosti i

predviđenih vrijednosti (Mishra, 2018). Srednju kvadratnu pogrešku izračunava se pomoću jednadžbe:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (26)$$

Korijen srednje kvadratne pogreške nastaje uzimanjem kvadratnog korijena srednje kvadratne pogreške te se koristi također kao standardna statistička mjera. Korijen srednje kvadratne pogreške izračunava se pomoću jednadžbe:

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (27)$$

4.1.2. Mjere kvalitete izvedbe za eksperimente klasifikacije

Zadatak klasifikacije podrazumijeva drugačija pravila za ocjenu pogreške. Prednost je da je broj izlaza diskretan, pa se može točno odrediti je li predviđanje uspješno ili nije na primjeru binarne klasifikacije (Bonnin, 2017). Time se u radu koriste mjere: točnost (engl. *Accuracy – ACC*), područje ispod krivulje osjetljivosti (engl. *Area Under the Receiver Operating Characteristic curve – AUC ROC*) i F1 mjera.

Točnost izračunava broj ispravnih predviđanja modela i to se izražava u postotku od 0 % do 100 %. Ovo je najčešće korišten mjerni pokazatelj kod rješavanja problema klasifikacije i izračunava se putem izraza:

$$\text{točnost} = \frac{\text{broj stvarno pozitivnih} + \text{broj stvarno negativnih}}{\text{ukupni broj}} \quad (28)$$

Područje ispod krivulje osjetljivosti koristi se kao mjera performanse binarnog klasifikatora. Što je veća površina ispod navedene krivulje, to je izvedba binarnog klasifikatora bolja.

Preciznost opisuje omjer ispravno predviđenih pozitivnih promatranja u odnosu na ukupan broj predviđenih pozitivnih promatranja.

$$Prec = \frac{TP}{TP + FP}$$

Opoziv opisuje omjer ispravno predviđenih pozitivnih promatranja u odnosu na sva predviđanja u promatranoj klasi.

$$Rec = \frac{TP}{TP + FN}$$

F1 mjera je harmonijska sredina preciznosti i opoziva te se često koristi mjera koja pokazuje pravu točnost modela. F1 mjera izračunava se pomoću jednadžbe:

$$F1 = 2 * \frac{preciznost * opoziv}{preciznost + opoziv} \quad (29)$$

4.1.3. Analiza statističkog značaja razlika

Kao pomoć pri interpretaciji rezultata dobivenih od modela predviđanja dodatno se provela analiza statističkog značaja razlika i analiza veličine učinka. Sukladno provođenju analize statističkog značaja razlika postavljaju se statističke hipoteze koje se potom testiraju.

Time postoje:

H_0 (nulta hipoteza) – nema statistički značajne razlike, odnosno dobivena razlika je rezultat slučajnih čimbenika.

H_1 (alternativna hipoteza) – postoji statistički značajna razlika te se time odbacuje nultu hipotezu.

Također osim postavljanja hipoteza odabire se nivo značajnosti alfa (α) koji pokazuje maksimalno dozvoljenu vjerojatnost greške na kojoj se odbacuje nultu hipotezu. Ako je vrijednost p provedenog statističkog testa $p < \alpha$, odbacuje se nultu hipotezu i prihvaća alternativnu.

Za testiranje statističkog značaja razlika nad rezultatima modela predviđanja ovog istraživanja koristi se parametarski statistički t-test (Studentov t-test) i neparametarski statistički

Wilcoxonov test dok se kod testiranja veličine efekta koristi Cohen's D test. Prije provođenja parametarskog statističkog t-testa nužno je provesti statističko testiranje normalnosti putem Kolmogorov-Smirnov testa. Ako se utvrdi statistički značajno odstupanje od normalne raspodjele, za taj uređeni par se ne provodi t-test.

4.2. Proces treniranja modela

Proces treniranja započinje odabirom jedne od kombinacija ulaznih podataka te njezinom podjelom na trening skup kojim se vrši treniranje modela, te na testni skup na kojem se u konačnici testiraju performanse modela. Drugi korak u procesu treniranja je postavljanje modela predviđanja i odabir njegovih parametara.

4.2.1. Parametri modela predviđanja

Pod izborom parametara modela predviđanja ubraja se određene varijable o kojima u konačnici ovise performanse modela u izvršavanju promatranih zadataka. Drugim riječima, riječ je o parametrima za koje treba pronaći najbolju moguću vrijednost izraženu kroz izabrane mjere kvalitete izvedbe na trening skupu podataka. U ovom istraživanju u te parametre ubraja se arhitekturu neuronske mreže koja se ogleda u vidu broja skrivenih slojeva i broja neurona unutar svakog sloja. Također, tu spadaju i vrijednosti vezane za algoritme optimizacije, kao što su vrijednosti stope učenja i broja epoha u procesu treniranja modela (Diaz *et al.*, 2017).

4.2.2. Odabir parametra modela predviđanja

Odabir najbolje postavke parametara rezultirat će u konačnici boljom izvedbom modela. Za svaki model posebno se pronalaze najbolje kombinacije parametara postavki što je dug i zamoran proces (Shevchuk, 2016). Postoji nekoliko definiranih načina i tehnika na koje se provodi ovaj proces gdje se navodi način nasumičnog odabira vrijednosti parametara uz korištenje tehnike pod nazivom višestruka unakrsna validacija (engl. *K-Fold Cross Validation*). Testiranja za odabir najboljih parametara postavke mreže provode se nad trening skupom podataka gdje se putem spomenute tehnike trening skup dodatno dijeli i na validacijski skup.

S ciljem pronalaska najboljih postavki parametara modela predviđanja korišten je nasumičan odabir u vidu da se svaki parametar testirao posebno. Kako se radi o modelima predviđanja koji ne obrađuju velike serije podataka, korištena je arhitektura povratne neuronske mreže koja se sastoji od ulaznog sloja, tri skrivena sloja i izlaznog sloja. Kroz testiranja za odabir najboljih postavki testirale su se vrijednosti prikazane u Tablici 4.

Tablica 4. Prikaz testiranih vrijednosti parametara modela predviđanja

Naziv parametra	Testirane vrijednosti
Broj neurona prvog skrivenog sloja	128, 256, 512
Broj neurona drugog skrivenog sloja	128, 256, 512
Broj neurona trećeg skrivenog sloja	32, 64
Vrijednost stope učenja	0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9
Dropout vrijednost	0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9
Broj epoha	1000, 2000, 3000, 4000, 5000

Na kraju provedenog nasumičnog odabira vrijednosti parametara uz korištenje tehnike višestruke unakrsne validacije u Tablici 5. se prikazuju konačni odabrani parametri. Model predviđanja koji se koristi za provođenje eksperimenata ovoga istraživanja s odabranim vrijednostima parametara prikazanih u Tablici 5. u daljnjem tekstu će se nazivati odabrani model predviđanja.

Tablica 5. Prikaz odabranih vrijednosti parametara modela postupkom višestruke unakrsne validacije za rješavanje regresijskog zadatka

Naziv parametra	Odabrane vrijednosti za regresijski model	Odabrane vrijednosti za klasifikacijski model
Broj neurona prvog skrivenog sloja	512	512
Broj neurona drugog skrivenog sloja	512	512
Broj neurona trećeg skrivenog sloja	32	32
Vrijednost stope učenja	0.3	0.3
Dropout vrijednost	0.3	0.3
Broj epoha	1000	3000

4.3. Rezultati eksperimenata

Uspješno provedenim eksperimentom došlo se do rezultata sukladno korištenim modelima za rješavanje regresijskog i klasifikacijskog zadatka. Prilikom prikaza rezultata koriste se ranije spomenuti mjerni pokazatelji u podjelama po ulaznim kombinacijama korištenim pri provođenju eksperimenta.

Za potrebe ovog rada uzele su se u razmatranje samo određene mjere nad kojima se dodatno radio i statistički značaj razlika kako bi samim time rad ostao u razumnim granicama opsega za diplomski rad.

4.3.1. Rezultati regresijskog modela

U tablicama 6. i 8. prikazani su rezultati odabranog regresijskog modela provedenog nad burzovnim indeksima Dow Jones i NASDAQ, te su u tablicama 7. i 9. prikazani rezultati analize statističkog značaja razlika.

U Tablici 7. vidi se da odabrani model kod burzovnog indeksa Dow Jones ostvaruje svoje najbolje rezultate u vidu korištenih mjera kvalitete izvedbe za eksperimente regresije prilikom četvrte ulazne kombinacije (K4), odnosno korištenjem zaključne vrijednosti i sentiment-podataka preuzetih s FinSentS portala. Četvrta ulazna kombinacija (K4) postiže najbolje prosječne vrijednosti korištenih mjera kvalitete izvedbe za eksperimente regresije, odnosno apsolutne pogreške (MAE) i korijena srednje prosječne kvadratne pogreške (RMSE).

Tablica 6. Rezultati odabranog regresijskog modela na skaliranom testnom skupu podataka za burzovni indeks Dow Jones

Kombinacija	Prosječna apsolutna pogreška (MAE)	Korijen srednje prosječne kvadratne pogreške (RMSE)
K1	0.01007	0.01290
K2	0.00850	0.01199
K3	0.00869	0.01163
K4	0.00788	0.01093
K5	0.00855	0.01198

Nakon dobivenih rezultata putem korištenih mjera kvalitete izvedbe za eksperimente regresije provodi se analizu statističkog značaja razlika gdje se kombinaciju koja je ostvarila najbolje rezultate u ovome slučaju K4 unosi kao jedan od parova s ostalim kombinacijama. Rezultati u Tablici 7. pokazuju da K4 s K1 i K3 pokazuje značajne statističke razlike dok s K2 i K5 ne pokazuje značajne statističke razlike. Također treba uzeti u obzir vrlo malu veličinu učinka kod svih uređenih parova što pokazuje da postoji veliko prekrivanje distribucija.

Tablica 7. Rezultati analize statističkog značaja razlika za mjeru prosječne apsolutne pogreške između odabranih uređenih parova za odabrani regresijski model burzovnog indeksa Dow Jones

Uređeni parovi	T-test			Cohen's D vrijednost	Wilcoxonov test		
	alpha	p-vrijednost	H ₀		alpha	p-vrijednost	H ₀
K4 s K1	0.01	< 0.001	odbaci	0.022	0.01	< 0.001	odbaci
K4 s K2	0.01	0.004	odbaci	-0.004	0.01	0.030	neuspjelo odbacivanje
K4 s K3	0.01	< 0.001	odbaci	0.011	0.01	<0.001	odbaci
K4 s K5	0.01	0.078	neuspjelo odbacivanje	0.002	0.01	0.979	neuspjelo odbacivanje

U Tablici 8. prikazani su rezultati odabranog regresijskog modela u vidu korištenih mjera kvalitete izvedbe za eksperimente regresije za burzovni indeks NASDAQ. Kao kod prethodnih rezultata za burzovni indeks Dow Jones, vidi se da su najbolje prosječne vrijednosti korištenih mjera kvalitete izvedbe za eksperimente regresije otvarani prilikom četvrte ulazne kombinacije (K4) kako kod rezultata prosječne apsolutne pogreške (MAE) tako i za korijen srednje kvadratne pogreške (RMSE).

Tablica 8. Rezultati odabranog regresijskog modela na skaliranom testnom skupu podataka za burzovni indeks NASDAQ

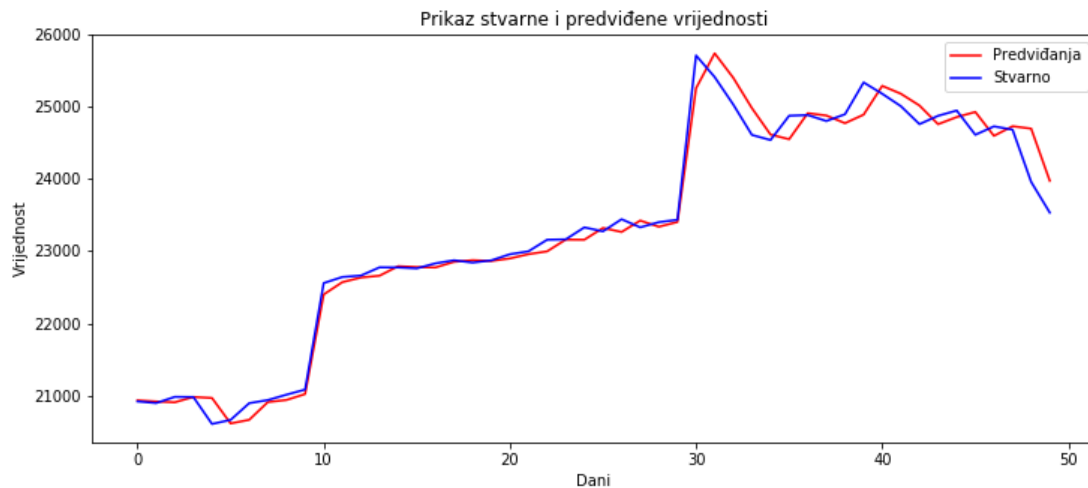
Kombinacija	Prosječna apsolutna pogreška (MAE)	Korijen srednje kvadratne pogreške (RMSE)
K1	0.00993	0.01276
K2	0.00869	0.01156
K3	0.00776	0.01075
K4	0.00720	0.01008
K5	0.00813	0.01133

Ponovno se provela analiza statističkog značaja razlika za burzovni indeks NASDAQ te se putem rezultata prikazanih u Tablici 9. utvrđuje da između svih uređenih parova postoji značajna statistička razlika s također malom veličinom učinka.

Tablica 9. Rezultati analize statističkog značaja razlika između odabranih uređenih parova za mjeru prosječne apsolutne pogreške za odabrani regresijski model burzovnog indeksa NASDAQ

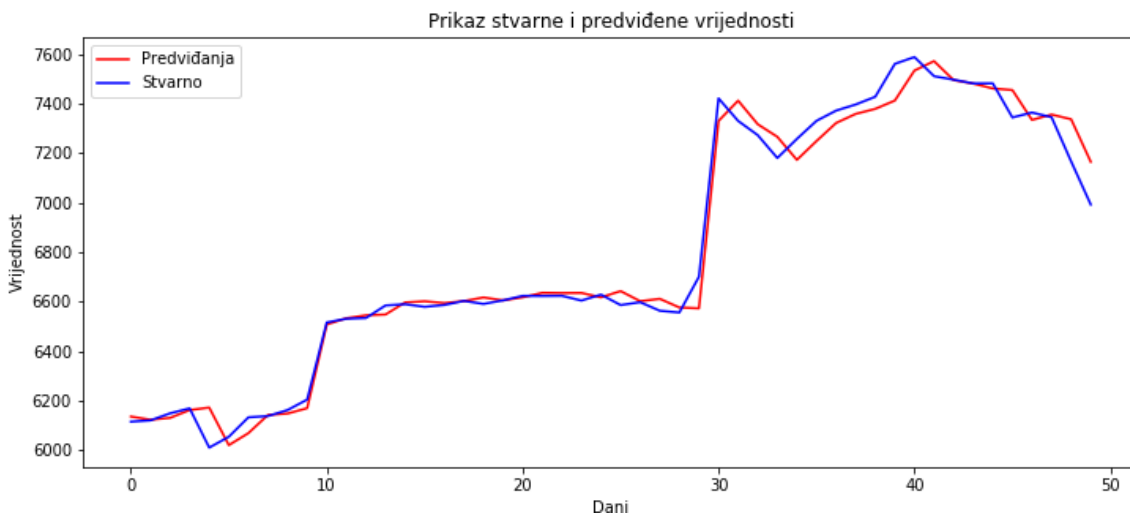
Uređeni parovi	T-test			Cohen's D vrijednost	Wilcoxonov test		
	alpha	p-vrijednost	H ₀		alpha	p-vrijednost	H ₀
K4 s K1	0.01	< 0.001	odbaci	0.028	0.01	< 0.001	odbaci
K4 s K2	0.01	< 0.001	odbaci	0.010	0.01	< 0.001	odbaci
K4 s K3	0.01	< 0.001	odbaci	-0.011	0.01	< 0.001	odbaci
K4 s K5	0.01	-	-	0.218	0.01	< 0.001	odbaci

Na slikama 16. i 17. grafički se prikazuju stvarne i predviđene vrijednosti dobivene korištenjem četvrte kombinacije koja je pokazala najbolje rezultate putem korištenih mjera kvalitete izvedbe za eksperimente regresije.



Slika 16. Prikaz stvarne i predviđene vrijednosti primjenom odabranog regresijskog modela za burzovni indeks Dow Jones

Izvor: izradio autor (2018)



Slika 17. Prikaz stvarne i predviđene vrijednosti primjenom odabranog regresijskog modela za burzovni indeks NASDAQ

Izvor: izradio autor (2018)

Sličnim istraživanjem Weng (2017) dolazi do gotovo identičnih rezultata koristeći model umjetnih neuronskih mreža. Za ulazne varijable koristio se u osnovi tehničkim indikatorima i FinSentS sentiment-podatcima čime je ostvario prosječnu apsolutnu pogrešku (MAE) u iznosu od 0.00820 i korijen srednje kvadratne pogreške (RMSE) u iznosu od 0.01130.

4.3.2. Rezultati klasifikacijskog modela

Provođenjem eksperimenta s odabranim klasifikacijskim modelom dobivaju se rezultati prikazani u tablicama 10. i 12., te su u tablicama 11. i 13. prikazani rezultati analize statističkog značaja razlika. Stavka unutar rezultata na koju će se obratiti posebnu pažnju jest točnost modela koja u konačnici govori koliko točno je model predvidio sutrašnje kretanje trenda.

Kao i kod odabranog regresijskog modela najbolji rezultati za burzovni indeks Dow Jones ostvareni su prilikom korištenja četvrte ulazne kombinacije (K4) u vidu korištenih mjera kvalitete izvedbe za eksperimente klasifikacije. Četvrta ulazna kombinacija (K4) ostvaruje najbolje prosječne vrijednosti korištenih mjera kvalitete izvedbe za eksperimente klasifikacije točnosti, AUC i F1 mjere. Posebno treba istaknuti rezultat točnosti modela koji iznosi visokih 60 % kao i rezultat F1 mjere koji iznosi također visokih 65 % za ovaj tip problema.

Tablica 10. Rezultati klasifikacijskog modela na skaliranom testnom skupu podataka za burzovni indeks Dow Jones

Kombinacija	Točnost	AUC	F1 mjera
K1	0.545	0.530	0.630
K2	0.541	0.527	0.623
K3	0.549	0.541	0.603
K4	0.605	0.598	0.651
K5	0.556	0.545	0.624

Provedenom analizom statističkog značaja razlika gdje se četvrta kombinacija (K4) unosi kao jedan od parova s ostalim kombinacijama, dobiveni rezultati prikazani su u Tablici 11. Rezultati u svim slučajevima pokazuju da nema statistički značajnih razlika između uređenih parova uz malu veličinu učinka. Ravnopravno se mogu koristiti svi testirani algoritmi vezano za mjeru točnosti za koju je testiran statistički značaj razlika.

Tablica 11. Rezultati analize statističkog značaja razlika između odabranih uređenih parova točnosti za odabrani klasifikacijski model burzovnog indeksa Dow Jones

Uređeni parovi	T-test			Cohen's D vrijednost	Wilcoxonov test		
	alpha	p-vrijednost	H ₀		alpha	p-vrijednost	H ₀
K4 s K1	0.01	0.015	neuspjelo odbacivanje	-0.204	0.01	0.015	neuspjelo odbacivanje
K4 s K2	0.01	0.039	neuspjelo odbacivanje	-0.179	0.01	0.039	neuspjelo odbacivanje
K4 s K3	0.01	0.927	neuspjelo odbacivanje	-0.007	0.01	0.927	neuspjelo odbacivanje
K4 s K5	0.01	0.249	neuspjelo odbacivanje	-0.007	0.01	0.248	neuspjelo odbacivanje

Rezultati odabranog klasifikacijskog modela za burzovni indeks NASDAQ prikazani su u Tablici 12. Iz tablice se vidi da su opet najbolji rezultati mjera kvalitete izvedbe za eksperimente klasifikacije ostvareni prilikom četvrte ulazne kombinacije (K4). Četvrta ulazna kombinacija (K4) ostvaruje najbolje prosječne vrijednosti korištenih mjera kvalitete izvedbe za eksperimente klasifikacije točnosti, AUC i F1 mjere.

Tablica 12. Rezultati klasifikacijskog modela na skaliranom testnom skupu podataka za burzovni indeks NASDAQ

Kombinacija	Točnost	AUC	F1 mjera
K1	0.560	0.501	0.713
K2	0.541	0.522	0.623
K3	0.556	0.520	0.672
K4	0.590	0.542	0.717
K5	0.568	0.539	0.665

Analizom statističkog značaja razlika između četvrte kombinacije koja je ostvarila najbolje rezultate mjera kvalitete izvedbe za eksperimente klasifikacije (istaknute vrijednosti u Tablici 12.) i ostalih kombinacija dobiva se da između njih postoje statistički značajne razlike uz srednju veličinu učinka. Za NASDAQ indeks K4 i RNN postiže najbolje rezultate s obzirom na testiranu mjeru točnosti (engl. *Accuracy – ACC*) što je vidljivo iz Tablice 13.

Tablica 13. Rezultati analize statističkog značaja razlika između odabranih uređenih parova za mjeru točnosti za odabrani klasifikacijski model burzovnog indeksa NASDAQ

Uređeni parovi	T-test			Cohen's D vrijednost	Wilcoxonov test		
	alpha	p-vrijednost	H ₀		alpha	p-vrijednost	H ₀
K4 s K1	0.01	0.0004	odbaci	-0.316	0.01	0.0005	odbaci
K4 s K2	0.01	-	-	0.564	0.01	< 0.001	odbaci
K4 s K3	0.01	0.0032	odbaci	0.255	0.01	0.0034	odbaci
K4 s K5	0.01	-	-	0.405	0.01	0.0034	odbaci

Postoje brojna istraživanja koja su zasnovana na sličnom cilju predviđanja kretanja trenda. Kad je riječ o korištenju samo tehničkih indikatora, Nelson, Pereira i De Oliveira (2017) ostvaruju rezultat točnosti modela u rasponu od 53,3 % do 55,9 %, dok Weng (2017) ostvaruje visokih 61,6 %. Wenga (2017) u svome istraživanju koristi umjetne neuronske mreže, stroj s potpornim vektorima i stabla odluke gdje navodi da su najbolji rezultati postignuti s metodom stroja s potpornim vektorima stoga su samo oni i prikazani.

Prema sličnom provedenom istraživanju, Atkins, Niranjana i Gerding (2018) testiraju također utjecaj Reuters novosti (Reuters, 2018) s burzovnim indeksima Dow Jones i NASDAQ. Ostvaruju rezultate od 56,6 % točnosti za predviđanje klasifikacijom trenda za Dow Jones, te 61,5 % za NASDAQ burzovni indeks. Kod ovog istraživanja autori su provodili predviđanja za svaki sljedeći sat te su dodatno koristili posebne metode za smanjenje dimenzija ulaznih značajki.

Najznačajniji rezultati kad je riječ o provođenju sentiment-analize za klasifikaciju trenda iznose visokih 73,55 % primjenom povratnih neuronskih mreža i novosti financijskih medija na području Kine (Liu *et al.*, 2017). Dodatno visokih 85,5 % ostvaruje Weng (2017) korištenjem posebno generiranih ulaznih značajki koje opisuje u radu, a one se temelje na statističkim pokazateljima preuzetim s Wikipedije i Googleove tražilice.

5. ZAKLJUČAK

Istraživanjem u ovom radu razvili su se modeli predviđanja putem kojih se može s određenom preciznošću predvidjeti kretanja na tržištu kapitala. Cilj je bio pokazati kako različiti pristupi u vidu tehničke, fundamentalne i sentiment-analize mogu pridonijeti boljim rezultatima pri korištenju povratnih neuronskih mreža koje su odabrane kao model predviđanja ovog rada.

Za izbor metode strojnog učenja koja se koristi u modelima predviđanja odabrane su umjetne neuronske mreže, točnije varijanta povratnih neuronskih mreža s ćelijama s dugoročnom memorijom. Ta metoda strojnog učenja pokazala se izrazito uspješnom u predviđanjima financijski vremenskih nizova te je kao takva korištena u brojnim relevantnim istraživanjima. Za ostvarenje ovoga cilja kreirane su razne kombinacije ulaznih podataka koje obuhvaćaju navedene analize te su korištene u kombinaciji s RNN-om. Modeli predviđanja podijeljeni su na rješavanje dvaju zadataka strojnog učenja, a to su regresija i klasifikacija.

Dobivenim rezultatima pokazalo se da se koristeći metode strojnog učenja može poboljšati predviđanja kretanja na tržištu kapitala i time pomoći trgovcima i investitorima u donošenju odluka. Najbolje prosječne vrijednosti korištenih mjera kvalitete izvedbe u predviđanju vrijednosti burzovnog indeksa i predviđanja kretanja trenda ostvareni su prilikom korištenja kombinacijom sa zaključnom vrijednosti burzovnog indeksa i FinSentS sentiment-podacima. Posebno se može istaknuti rezultat od 60,5 % točnosti kod odabranog klasifikacijskog modela nad burzovnim indeksom Dow Jones, te prosječna apsolutna pogreška od 0.00720 kod odabranog regresijskog modela nad burzovnim indeksom NASDAQ.

Ovim istraživanjem otvara se prostor za brojna poboljšanja modela predviđanja. Jedna od ideja za buduća istraživanja je razvoj automatiziranih sustava za trgovanje na elektroničkim burzama koristeći pri tome tehnike strojnog učenja. S trenutnim ostvarenim rezultatima odabranih modela predviđanja mogu se implementirati brojne strategije trgovanja u vidu generiranja kupovnih i prodajnih signala ne bi li se time trgovcima olakšalo proces donošenja odluke. Strojno učenje s primjenom u području ekonometrije otvara sasvim novu domenu koja je sa svojim razvojem tek započela.

POPIS LITERATURE

- Alpaydin, E. (2009) *Introduction to machine learning*, MIT press.
- Alpha Vantage (2018) *Alpha Vantage*. <raspoloživo na: <https://www.alphavantage.co/>>, pristupljeno [10/06/2018].
- Arlt, J. i Arltová, M. (2001) 'Financial Time Series and Their Features', *Acta oeconomica pragensia VŠE Praha*, 9(4), pp. 7–20.
- Atkins, A., Niranjan, M. i Gerding, E. (2018) 'Financial news predicts stock market volatility better than close price', *The Journal of Finance and Data Science*, Elsevier Ltd.
- Barbić, T. (2010) 'Pregled razvoja hipoteze efikasnog tržišta', *Privredna kretanja i ekonomska politika*, 124, pp. 29–61.
- Baresa, S., Bogdan, S. i Ivanovic, Z. (2013) 'Strategy of Stock Valuation By Fundamental Analysis', *UTMS Journal of Economics*, 4(1), pp. 45–51.
- Birău, F. R. (2012) 'The Impact of Behavioral Finance on Stock Markets', *Economy Series*, 78421(3), pp. 45–50.
- Blum, A. (2007) 'Machine Learning Theory'.
- Bonnin, R. (2017) *Machine Learning for Developers*, Packt Publishing.
- Brajković, A. i Peša, A. (2015) 'Bihevioralne financije i teorija „Crnog labuda“ Behavioral Finance and "Black Swan" Theory', *Oeconomica Jadertina*, pp. 65–93.
- Brownlee, J. (2016a) *How to Grid Search Hyperparameters for Deep Learning Models in Python With Keras*. <raspoloživo na: <https://machinelearningmastery.com/grid-search-hyperparameters-deep-learning-models-python-keras/>>, pristupljeno [02/07/2018].
- Brownlee, J. (2016b) *Time Series Prediction with LSTM Recurrent Neural Networks in Python with Keras*. <raspoloživo na: <https://machinelearningmastery.com/time-series-prediction-lstm-recurrent-neural-networks-python-keras/>>.
- Calderon, P. (2017) *VADER Sentiment Analysis Explained*. <raspoloživo na: <http://datameetsmedia.com/vader-sentiment-analysis-explained/>>, pristupljeno [29/06/2018].
- Colby, R. W. (2003) *The Encyclopedia of Technical Market Indicators*, McGraw-Hill.

- Dan-Ching, L. (2017) *A Practical Guide to ReLU*. <raspoloživo na: <https://medium.com/tiny-mind/a-practical-guide-to-relu-b83ca804f1f7>>, pristupljeno [27/06/2018].
- Darskuvienė, V. (2010) *Financial Markets, Leonardo da Vinci program project*.
- Delbelo Bašić, B., Čupić, M. i Šnajder, J. (2008) *Umjetne neuronske mreže*, Fakultet elektrotehnike i računarstva.
- Diaz, G., Fokoue, A., Nannicini, G. i Samulowitz, H. (2017) 'An effective algorithm for hyperparameter optimization of neural networks', *IBM Journal of Research and Development*, 61(4), pp. 1–11.
- Donges, N. (2018) *Recurrent Neural Networks and LSTM*. <raspoloživo na: <https://towardsdatascience.com/recurrent-neural-networks-and-lstm-4b601dd822a5>>, pristupljeno [28/05/2018].
- Fama, E. F. (1970) 'Efficient Capital Markets - A Review of Theory and Empirical Work.pdf', *Journal of Finance*, pp. 383–417.
- FinSentS (2018) *FinSentS*. <raspoloživo na: <http://landing.finsents.com/>>, pristupljeno [10/06/2018].
- Golik, P., Doetsch, P. i Ney, H. (2013) 'Cross-entropy vs. Squared error training: A theoretical and experimental comparison', *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*, 2(2), pp. 1756–1760.
- Goodfellow, I., Bengio, Y. and Courville, A. (2016) 'Deep learning', *Healthcare Informatics Research*, 22(4), pp. 351–354.
- Hansson, M. (2017) 'On stock return prediction with LSTM networks'.
- Hastie, T., Tibshirani, R. i Friedman, J. (2001) 'The Elements of Statistical Learning', *The Mathematical Intelligencer*, 27(2), pp. 83–85.
- Haykin, S. (2009) *Neural Networks and Learning Machines*, Pearson Prentice Hall.
- Hochreiter, S. i Schmidhuber, J. (1997) 'Long Short-Term Memory', *Neural Computation*, 9(8), pp. 1735–1780.
- Hutto, C. J. i Gilbert, E. (2014) 'Vader: A parsimonious rule-based model for sentiment analysis of social media text', in *Eighth International AAAI Conference on Weblogs and ...*,

pp. 216–225.

InfoTrie (n.d.) 'InfoTrie FinSentS Static Data User Guide'. <raspoloživo na: <http://landing.finsents.com/>>.,

Investopedia (2018) *The Top Sites for the Latest Stock Market News | Investopedia*. <raspoloživo na: <https://www.investopedia.com/articles/investing/112514/top-sites-latest-stock-market-news.asp>>, pristupljeno [02/06/2018].

Investopedia (n.d.) *Nasdaq Composite Index*. <raspoloživo na: <https://www.investopedia.com/terms/n/nasdaqcompositeindex.asp>>, pristupljeno [28/05/2018].

Investopedia (n.d.) *Technical Indicator*. <raspoloživo na: <https://www.investopedia.com/terms/t/technicalindicator.asp>>, pristupljeno [10/06/2018].

Kapur, R. i Khazan, L. (2016) *Neural Networks & The Backpropagation Algorithm, Explained*. <raspoloživo na: <https://ayearofai.com/rohan-lenny-1-neural-networks-the-backpropagation-algorithm-explained-abf4609d4f9d>>, pristupljeno [09/06/2018].

Keras (2018) *Keras*. <raspoloživo na: <https://keras.io/>>, pristupljeno [02/07/2018].

Krein, D. i W. Smith, J. (2014) *The NASDAQ Composite Index*. <raspoloživo na: <https://blogs.wsj.com/public/resources/documents/nasdaqcompretperspective.pdf>>.,

Kriesel, D. (2005) *A Brief Introduction to Neural Networks*. <raspoloživo na: http://www.dkriesel.com/_media/science/neuronalenetze-en-zeta2-2col-dkrieselcom.pdf>.,

Krugman, P. (2009) *How Did Economists Get It So Wrong? - The New York Times*. <raspoloživo na: <https://www.nytimes.com/2009/09/06/magazine/06Economic-t.html>>, pristupljeno [02/07/2018].

Liu, Y., Zengchang, Q., Pengyu, L. i Tao, W. (2017) 'Stock volatility prediction using recurrent neural networks with sentiment analysis', in *International Conference on Industrial, Engineering and Other Applications of Applied Intelligent Systems*, pp. 31–40.

Maas, A. L., Hannun, A. Y. i Ng, A. Y. (2013) 'Rectifier Nonlinearities Improve Neural Network Acoustic Models', in *Proceedings of the 30th International Conference on Machine Learning*, p. 6. <raspoloživo na: https://web.stanford.edu/~awni/papers/relu_hybrid_icml2013_final.pdf>.,

Malkiel, B. G. (2003) 'The Efficient Market Hypothesis and Its Critics', *Journal of Economic Perspectives*, 17(1), pp. 59–82.

Marsland, S. (2009) *An Algorithmic Perspective*, Chapman & Hall/CRC.

McCulloch, W. S. i Pitts, W. (1943) 'A logical calculus of the ideas immanent in nervous activity', *The Bulletin of Mathematical Biophysics*, 5(4), pp. 115–133.

Mishra, A. (2018) *Metrics to Evaluate your Machine Learning Algorithm*. <raspoloživo na: <https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-f10ba6e38234>>, pristupljeno [06/06/2018].

Mohammed, M., Khan, M. B. i Mohammed Bashier, E. B. (2006) *Machine learning: algorithms and applications*, CRC Press.

Munoz, A. (2018) *Stock-Prediction*. <raspoloživo na: <https://github.com/munozalexander/Stock-Prediction>>, pristupljeno [02/07/2018].

Murphy, J. J. (1999) *Technical Analysis in Financial Markets*, New York Institute of Finance. doi: 10.2139/ssrn.566882.

Nargesian, F., Samulowitz, H., Khurana, U., Khalil, E. B. i Turaga, D. (2017) 'Learning Feature Engineering for Classification', in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI*.

Nelson, D. M. Q., Pereira, A. C. M. i De Oliveira, R. A. (2017) 'Stock market's price movement prediction with LSTM neural networks', in *Proceedings of the International Joint Conference on Neural Networks*. doi: 10.1109/IJCNN.2017.7966019.

Nielsen, M. A. (2015) *Neural Networks and Deep Learning*, Determination Press.

Nielsen, M. A. (2017) 'Neural Networks and Deep Learning'. <raspoloživo na: <http://neuralnetworksanddeeplearning.com/chap1.html>>.,

Olah, C. (2015) *Understanding LSTM Networks*. <raspoloživo na: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>>, pristupljeno [29/05/2018].

P. Murphy, K. (2012) *Machine Learning: A Probabilistic Perspective*, The MIT Press.

Pandas (2018) *Pandas*. <raspoloživo na: <https://pandas.pydata.org/>>, pristupljeno [02/07/2018].

Patel, J., Shah, S., Thakkar, P. i Kotecha, K. (2015a) 'Predicting stock and stock price index

movement using Trend Deterministic Data Preparation and machine learning techniques’, *Expert Systems with Applications*, Elsevier Ltd, 42(1), pp. 259–268.

Patel, J., Shah, S., Thakkar, P. i Kotecha, K. (2015b) ‘Predicting stock market index using fusion of machine learning techniques’, *Expert Systems with Applications*, Elsevier Ltd, 42(4), pp. 2162–2172.

Pawar, A. B., Jawale, M. A. i Kyatanavar, D. N. (2016) *Fundamentals of Sentiment Analysis Concepts and Methodology*, *Studies in Computational Intelligence*.

Person, J. L. (2004) *Complete Guide to Technical Trading Tactics*, Wiley Trading.

Petrusheva, N. i Jordanoski, I. (2016) ‘Comparative Analysis Between the Fundamental and Technical Analysis of Stocks’, *Journal of Process Management. New Technologies*, 4(2), pp. 26–31.

Pimprikar, R., Ramachandran, S. i Senthilkumar, K. (2017) ‘Use of machine learning algorithms and twitter sentiment analysis for stock market prediction’, *International Journal of Pure and Applied Mathematics*, 115(6), pp. 521–526.

Python (2018) *Python*. <raspoloživo na: <https://www.python.org/>>, pristupljeno [02/07/2018].

Ramachandran, P., Zoph, B. i Le, Q. V. (2017) ‘Searching for Activation Functions’, pp. 1–13.

Rechenthin, M. D. (2014) *Machine-learning classification techniques for the analysis and prediction of high-frequency stock direction*, *Iowa Research Online*, University of Iowa.

Reuters (2018) *Reuters*. <raspoloživo na: <https://www.reuters.com/>>, pristupljeno [10/06/2018].

Romeu, R. (2001) *Technical Analysis for Direct Access Trading: A Guide to Charts, Indicators & Other Indispensable Market Analysis Tools*, McGraw-Hill.

Schumaker, R. P. i Chen, H. (2009) ‘Textual analysis of stock market prediction using breaking financial news’, *ACM Transactions on Information Systems*, 27(2), pp. 1–19.

Shalev-Shwartz, S. i Ben-David, S. (2014) *Understanding machine learning: From theory to algorithms*, Cambridge University Press.

Shevchuk, Y. (2016) *Hyperparameter optimization for Neural Networks*. <raspoloživo na: http://neupy.com/2016/12/17/hyperparameter_optimization_for_neural_networks.html#hyper

parameter-optimization>, pristupljeno [07/06/2018].

Shoven, J. B. (2000) 'The Dow Jones Industrial Average : The Impact of Fixing Its Flaws', *Journal of Wealth Management*.

Silva, I. N., Hernane Spatti, D., Andrade Flauzino, R., Liboni, L. H. B. i dos Reis Alves, S. F. (2017) 'Artificial Neural Networks', in *IJCAI International Joint Conference on Artificial Intelligence*, Springer, p. 307.

Šlibar, E. (2009) *Metode izračuna burzovnih indeksa*.

Tandel, A. (2017) *Support Vector Machines — A Brief Overview – Towards Data Science*. <raspoloživo na: <https://towardsdatascience.com/support-vector-machines-a-brief-overview-37e018ae310f>>, pristupljeno [09/06/2018].

The Royal Society (2017) *Machine learning: the power and promise of computers that learn by example*. <raspoloživo na: royalsociety.org/machine-learning>.,.

Thomsett, M. C. (1998) 'Mastering Fundamental Analysis'.

Tsai, C. F. i Wang, S. P. (2009) 'Stock Price Forecasting by Hybrid Machine Learning Techniques', in *Proceedings of the International MultiConference of Engineers and Computer Scientists*.

Tsay, R. S. (2010) *Analysis of Financial Time Series: Third Edition*, John Wiley & Sons.

Weng, B. (2017) *Application of machine learning techniques for stock market prediction*, Graduate Faculty of Auburn University.

Widegren, P. (2017) *Deep learning-based forecasting of financial assets*, KTH Royal Institute Of Technology School Of Engineering Sciences.

Yang, L., Hanneke, S. i Carbonell, J. (2013) 'A Theory of Transfer Learning with Applications to Active Learning', *Machine learning*, pp. 1–28.

Zhang, C. i Ma, Y. (2012) *Ensemble Machine Learning*, Springer. doi: 10.1007/978-1-4419-9326-7.

POPIS SLIKA

Slika 1. Linearna regresija	6
Slika 2. Binarna klasifikacija	6
Slika 3. Simulirani podatci u ravnini, grupirani u tri klase	7
Slika 4. Biološki neuron.....	11
Slika 5. Model neurona.....	12
Slika 6. Sigmoid funkcija.....	14
Slika 7. Tanh funkcija.....	14
Slika 8. ReLU funkcija	15
Slika 9. Primjer višeslojne umjetne neuronske mreže	16
Slika 10. Gradijentni spust.....	17
Slika 11. Algoritam povratne propagacije	18
Slika 12. Kretanje burzovnog indeksa „Dow Jones“	20
Slika 13. Povratna neuronska mreža.....	26
Slika 14. Čelija s dugoročnom memorijom	26
Slika 15. Prikaz pomičnog prozora.....	37
Slika 16. Prikaz stvarne i predviđene vrijednosti primjenom odabranog regresijskog modela za burzovni indeks Dow Jones	46
Slika 17. Prikaz stvarne i predviđene vrijednosti primjenom odabranog regresijskog modela za burzovni indeks NASDAQ	46

POPIS TABLICA

Tablica 1. Prikaz vrijednosti preuzetih s mrežne baze podataka FinSentS za burzovni indeks Dow Jones	34
Tablica 2. Popis kombinacija ulaznih atributa za provođenje eksperimenata predviđanja vrijednosti burzovnih indeksa ili klase rasta/pada njihove vrijednosti	35
Tablica 3. Prikaz izvornih i skaliranih vrijednosti zaključne vrijednosti i vrijednosti preuzetih s mrežne baze podataka FinSentS za burzovni indeks Dow Jones	36
Tablica 4. Prikaz testiranih vrijednosti parametara modela predviđanja	42
Tablica 5. Prikaz odabranih vrijednosti parametara modela postupkom višestruke unakrsne validacije za rješavanje regresijskog zadatka	42
Tablica 6. Rezultati odabranog regresijskog modela na skaliranom testnom skupu podataka za burzovni indeks Dow Jones	43
Tablica 7. Rezultati analize statističkog značaja razlika za mjeru prosječne apsolutne pogreške između odabranih uređenih parova za odabrani regresijski model burzovnog indeksa Dow Jones	44
Tablica 8. Rezultati odabranog regresijskog modela na skaliranom testnom skupu podataka za burzovni indeks NASDAQ	45
Tablica 9. Rezultati analize statističkog značaja razlika između odabranih uređenih parova za mjeru prosječne apsolutne pogreške za odabrani regresijski model burzovnog indeksa NASDAQ	45
Tablica 10. Rezultati klasifikacijskog modela na skaliranom testnom skupu podataka za burzovni indeks Dow Jones.....	47
Tablica 11. Rezultati analize statističkog značaja razlika između odabranih uređenih parova točnosti za odabrani klasifikacijski model burzovnog indeksa Dow Jones	48
Tablica 12. Rezultati klasifikacijskog modela na skaliranom testnom skupu podataka za burzovni indeks NASDAQ.....	48
Tablica 13. Rezultati analize statističkog značaja razlika između odabranih uređenih parova za mjeru točnosti za odabrani klasifikacijski model burzovnog indeksa NASDAQ	49