

Univerzitet u Sarajevu  
Elektrotehnički fakultet  
Odsjek za Računarstvo i informatiku  
Dubinska analiza podataka  
Dokumentacija seminarskog rada

**Predicting Stock Market Movement  
Based on Twitter Data and News Articles  
Using Sentiment Analysis and Fuzzy Logic**

Džigal Džemil  
Hodo Midhat  
Piriya Ena

Sarajevo, juni 2021

## 1) Odabir problema za rješavanje, motivacija, specifikacija problema koji se rješava.

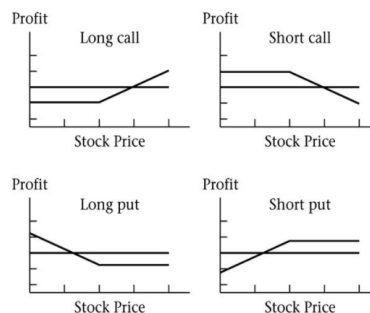
### Džemil:

Investiranje u dionice je dugogodišnja i temeljno ustanovljena struka koja zahtijeva godine studiranja i učenja na greškama kako bi osoba mogla biti iole uspješna u istoj. Generalno, dionice mogu da rastu i analogno tome, da opadaju. Najveći problem je predvidjeti u kojem smjeru će se kretati cijena dionice. Shodno tome, postoje razni mehanizmi “klađenja” za i protiv neke određene dionice. Ta “klađenja” se zovu stock calls i stock puts. Kupovanje stock calls je, pojednostavljeno, klađenje da će do isteka tog call-a cijena dionice porasti do i iznad cijene X. Kupovanje stock puts je, pojednostavljeno, klađenje da će do isteka tog put-a cijena dionice opasti do i ispod cijene X. Jasno je odakle dolazi faktor “klađenje”. Shodno tome, za maksimizaciju profita, potrebno je implementirati sistem za predlaganje koji korisniku implicitno predlaže da li treba kupiti short-term call ili put na određenu dionicu.



### Ena:

Jasno je iz prethodnog izlaganja da je vrijeme isticanja call i put poziva u uskoj vezi sa mogućim profitom nakon uspjeha call ili put-a. Što je vrijeme isticanja call-a i put-a kraće, to je rizik veći zbog prirodne fluktuacije svake dionice ali profiti mogu biti veliki. Analogno tome duga vremena trajanja rezultiraju skupljim call i put-ovima ali većom vjerovatnoćom za profit koji je, u ovom slučaju, manji nego za kraće call-ove i put-ove. Očigledno je da sistem koji predviđa smjer kretanja cijene dionice može primarno pomoći za kupovinu call-ova i put-ova sa kraćim vremenskim trajanjem.



### Midhat:

Jedna dionica koja objedinjuje stanje na svjetskoj berzi u odnosu na prvih 500 kompanija po cijeni njihovih dionica. Ovaj indeks, historijski, dobro opisuje kretanje čitave berze, primarno u SAD-u. Tako da ako možemo formirati neki model koji određuje short-term kretanje cijene ovog indeksa, možemo kupovati call-ove i put-ove za ovaj indeks. Shodno tome, potrebno je prikupiti

informacije koje generalno opisuju dionice. Jedan koristan izvor ovih informacija je Twitter. Pokazano je da twitter postovi su u korelaciji sa kretanjem cijena SNP500 indeksa.



Sve ukazuje na to da je moguće riješiti problem predviđanja kretanja cijena dionica na berzi tako što ćemo iskoristiti Twitter objave u sistemu kako bi se moglo odrediti kretanje berze, shodno tome i kretanje SNP500 indeksa. Dalje je jednostavno inferirati šta treba i kada treba kupiti. Ovakav sistem je više prikladan Sistemima podrške odlučivanju i možemo ga posmatrati kao kompromis između klasičnih ML problema čija rješenja su eksplicitna i opštih, sugestivnih problema ML-a čija rješenja su generalni prijedlozi šta treba uraditi (iz razloga što dobijamo sugestiju u kojem smjeru će se kretati cijene dionica odakle možemo inferirati koju odluku trebamo donijeti kao akter, odnosno investitor).

**Džemil, Ena, Midhat:**

Zaključujemo da bi ovaj sistem bio najviše koristan običnim investitorima, amaterima. Također, i veliki investitori u hedge fund kompanijama mogu koristiti sistem za svrhe poteza “velike kocke” koji se često dešavaju u ovim kompanijama, s obzirom na to da je uvijek u pitanju kompanija čiji kapital je reda milijardi USD.

## 2) Odabir naučnog rada koji je povezan sa temom i analiza istog.

**Džemil, Ena, Midhat:**

Rad koji smo odabrali za analizu i iz kojeg smo preuzeli mnoge ideje za implementaciju svojeg rada je Predicting Stock Market Movement Based on Twitter Data and News Articles Using Sentiment Analysis and Fuzzy Logic, Roshni et. al.

**Midhat:**

Uvod i Related Work:

U uvodu autori opisuju veliki uticaj socijalnih mreža na razne grane industrije. Jedna od tih grana je ulaganje u dionice. Autori ističu činjenicu da objave na socijalnim mrežama mogu formirati “kolektivno mišljenje” o stanju na berzi koje se može koristiti kako bi se kretanje iste moglo predviđati. Također spominju da kretanje Dow Jones indeksa može biti predviđeno sa tačnosti 87.6% koristeći samo Twitter podatke. Napominju da veliki uticaj na cijene dionica imaju poznate ličnosti, muzičari, poduzetnici, generalno bogati ljudi. Ovo ima smisla jer je poznato da ljudi sa uticajem mogu da sa tim svojim uticajem značajno utiču na opadanje odnosno rast cijena

određenih dionica i resursa. Ovo se očituje u posljednjem rastu i padu cijena kriptovaluta za čije haotično ponašanje se “glavnim krivcem” smatra poduzetnik Elon Musk.



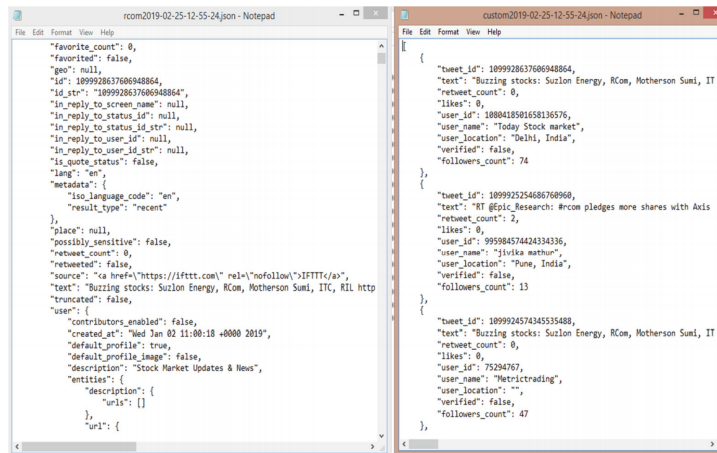
Dalje, autori opisuju sentimentalnu analizu u smislu rješavanja problema ternarne klasifikacije. Opisuju da ova analiza može biti iskorištena za formiranje robustnog ML modela sa većom tačnošću i najavljuju grafike performansi određenih modela koje su iskoristili. U dijelu Related Work opisuju razne prijašnje autore koji su iskoristili Random Forest klasifikatore, linearnu regresiju i SVM za postizanje sličnih rezultata, ali sa manjom tačnošću.

## Ena:

### Data Acquisition i Data Preprocessing:

U Data Acquisition dijelu autori opisuju način prikupljanja podataka koristeći Twitter-ove API-eve. Koristili su poznati paket tweepy i učitali su mnogo tweet-ova u JSON formatu. Ovim načinom su formirali bazu podataka koja predstavlja bazu znanja njihovog sistema. Naravno, ove podatke je potrebno obraditi. Autori su to opisali u dijelu Data Preprocessing u kojem navode da je potrebno ukloniti razne “problematične” simbole kao što su URL-ovi, EMOJI simboli i slični simboli koje smo i mi morali uklanjati pri radu sa tekstom na ovom predmetu. Dalje, tokenizirali

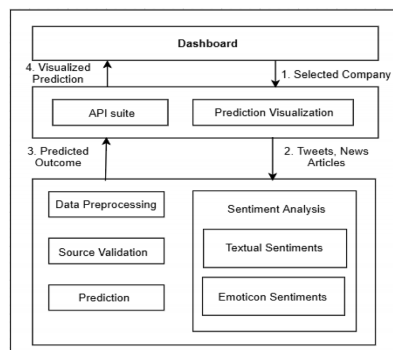
su tekst tweet-ova i, koristeći regularne ekspresije, izvršili su pattern matching kako bi uklonili hashtag-ove i spominjanje određenih korisnika, zamijenivši ih kao (#Nestle u Nestle i @DonaldTrump u user). Ovo ima za cilj maksimalno poopštenje konteksta teksta tweet-a koji se analizira.



## Džemil:

### Sentiment Analysis i Implementation:

Autori su izvršili sentimentalnu analizu koristeći NLTK biblioteku. Konkretno iskoristili su Naive Bayes Classifier unutar NLTK-a kako bi izvršili analizu. Dalje, izvršili su shallow parsing da odrede glavne dijelove rečenica kako bi se umanjio efekat outlier riječi. U dijelu preprocessing-a su uklonili određene emoji-e ali s obzirom na to da za ovu analizu treba i posljedica, odnosno labels za svaki tweet, određeni emoji-i se koriste kako bi se odredio sentiment tog tweet-a. Također, prioritet su imali tweet-ovi koji imaju veliki broj retweet-ova. Obrazložili su da tweet-ovi koji su više puta retweet-ovani imaju veći uticaj od drugih, što ima smisla jer se tim povećava doseg tog tweet-a pa i uticaj tog tweet-a na cijenu dionica, kao i tumačenje u kojem smjeru treba ići cijena dionice. Tweet-ovi korisnika sa manje od 100 follower-a se odbacuju jer nemaju dovoljno veliki doseg širokih narodnih masa. Predložili su izgled arhitekture sistema za konačnu kalkulaciju sentimenta odnosno odluke u kojem smjeru će se kretati cijena dionica na berzi.



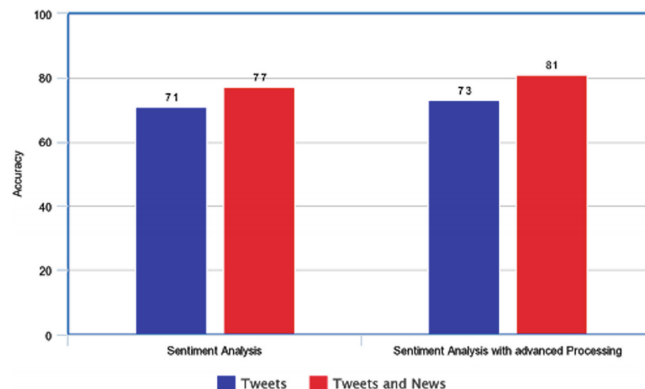
Što se tiče same implementacije, autori opisuju korake njihovog pristupa. Počinju sa akvizicijom podataka kako je opisano u dijelu prije ovog. Preprocesiraju te podatke, također kako je opisano u dijelu prije ovog. Dalje, izvršavaju sentimentalnu analizu koristeći tehniku sa emoji-ima koji

opisuju sentiment tog tweet-a i formiraju inicijalni zaključak o sentimentu tweet-a. Ovdje dolaze do zaključka da samo jedan izvor neće omogućiti maksimalnu tačnost. Zaključuju da trebaju uključiti i podatke iz novina kako bi povećali tačnost. Ovim formiraju hibridni NLP sistem sa dvije baze znanja (tweet-ovi i novinski članci) i koriste ih kako bi zaključili sentiment ovih podataka.

**Džemil, Midhat i Ena:**

Results and Discussion:

Kako je opisano u prošlom odjeljku, rezultati ukazuju na to da uključivanje novinskih članaka znatno povećava tačnost predikcije sentimentalne analize.



Ukazuju na značajno povećavanje tačnosti, u prosjeku oko 7-10%, u ovisnosti od toga koji podaci se koriste za analizu. U dijelu diskusije (zaključak i budući rad) autori ukazuju na moć RNN-ova za povećavanje tačnosti. Shodno njihovoj osobini da razmatraju eksponencijalno više različitih scenarija u odnosu na klasične ML algoritme i metode. Ovu metodu ostavljaju za dalji rad pri čemu sugerišu da je moguće postići tačnosti oko 81%. Ne opisuju kako se dolazi do te tačnosti, ali uz pomoć RNN-ova nije teško zaključiti kako (uz dovoljno velike mreže). Također, predlažu fuzzy sisteme regulacije za adaptaciju parametara RNN-ova kako bi se povećala tačnost istih. Također, zbog korištenja tweepy paketa, predlažu da uz dodatnu promjenu data acquisition koraka je moguće predvidjeti kretanje cijene dionica konkretnih kompanija, a ne samo SNP 500 indeksa i sličnih indeksa metrika svjetske berze.

- 3) **Prikupljanje i odabir podataka za sistem specificiran u 1. Priprema i preprocesiranje podataka (priložiti dokumentaciju i kodove za ovaj korak). Odabir algoritama mašinskog i dubokog učenja za rješavanje problema. Tabela prikaz algoritama sa naznakom ko je od članova tima zadužen za detaljniju analizu i implementaciju pojedinog algoritma.**

Za podatke je korišten sljedeći dataset:  [Stock Market Tweet | Sentiment Analysis lexicon | Kaggle](#)

Što se tiče preprocesiranja podataka korištena su sljedeća pomagala:

```
import warnings
warnings.simplefilter('ignore')
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns
sns.set()

from pprint import pprint
from datetime import datetime
import collections
import re

import nltk
nltk.download('stopwords')
nltk.download('punkt')
nltk.download('wordnet')
from nltk.corpus import wordnet, stopwords
from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize

from wordcloud import WordCloud

[nltk_data] Downloading package stopwords to
[nltk_data]   /home/chilichipsba/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to
[nltk_data]   /home/chilichipsba/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package wordnet to
[nltk_data]   /home/chilichipsba/nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
```

```
import re

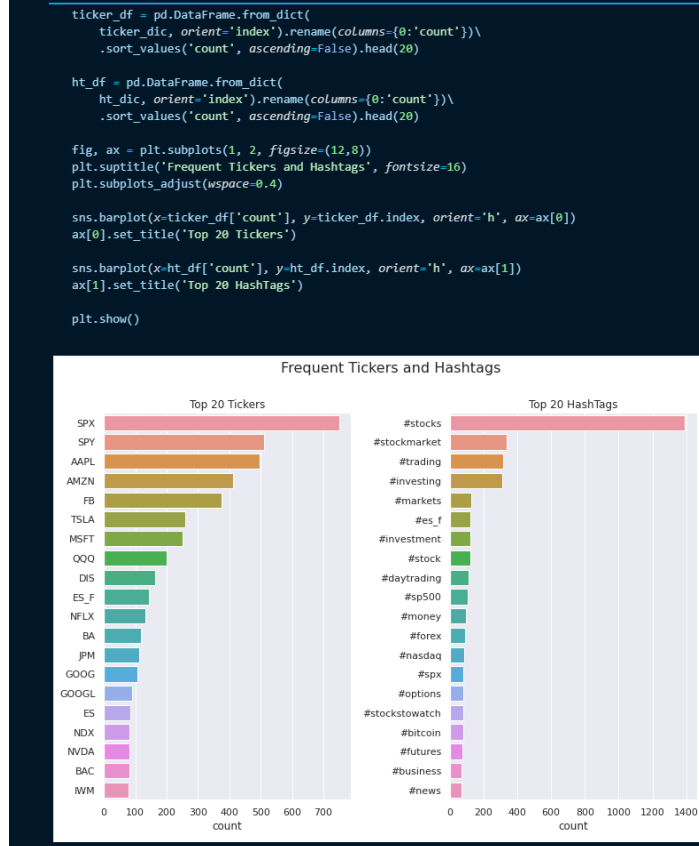
ticker_pattern = re.compile(r'^\s*[A-Z]+/\s*$ES_F$')
ht_pattern = re.compile(r'^\s*[w]+$')

ticker_dic = collections.defaultdict(int)
ht_dic = collections.defaultdict(int)

for text in data['text']:
    for word in text.split():
        if ticker_pattern.fullmatch(word) is not None:
            ticker_dic[word[1:]] += 1

        word = word.lower()
        if ht_pattern.fullmatch(word) is not None:
            ht_dic[word] += 1
```

[7]   M4



```
def decontracted(phrase):
    # ...
```

```

charency = re.compile('[-a-zA-Z]')
handle_pattern = re.compile('[-\w]')
emoji_pattern = re.compile('')

    "U+00000000-U+00000064" # emoticons
    "U+00000065-U+0000006F" # symbols & pictographs
    "U+00000070-U+0000007F" # transport & map symbols
    "U+00000100-U+0000013F" # flags (105)
    "U+00000140-U+0000017F"
    "U+00000180-U+000001FF"
    "U+00000200-U+000002FF"
    "U+00000300-U+000003FF"
    "U+00000400-U+000004FF"
    "U+00000500-U+000005FF"
    "U+00000600-U+000006FF"
    "U+00000700-U+000007FF"
    "U+00000800-U+000008FF"
    "U+00000900-U+000009FF"
    "U+00000A00-U+00000AFF"
    "U+00000B00-U+00000BFF"
    "U+00000C00-U+00000CFF"
    "U+00000D00-U+00000DFF"
    "U+00000E00-U+00000EFF"
    "U+00000F00-U+00000FFF"
    "U+00001000-U+000010FF"
    "U+00001100-U+000011FF"
    "U+00001200-U+000012FF"
    "U+00001300-U+000013FF"
    "U+00001400-U+000014FF"
    "U+00001500-U+000015FF"
    "U+00001600-U+000016FF"
    "U+00001700-U+000017FF"
    "U+00001800-U+000018FF"
    "U+00001900-U+000019FF"
    "U+00001A00-U+00001AFF"
    "U+00001B00-U+00001BFF"
    "U+00001C00-U+00001CFF"
    "U+00001D00-U+00001DFF"
    "U+00001E00-U+00001EFF"
    "U+00001F00-U+00001FFF"
    "U+00002000-U+000020FF"
    "U+00002100-U+000021FF"
    "U+00002200-U+000022FF"
    "U+00002300-U+000023FF"
    "U+00002400-U+000024FF"
    "U+00002500-U+000025FF"
    "U+00002600-U+000026FF"
    "U+00002700-U+000027FF"
    "U+00002800-U+000028FF"
    "U+00002900-U+000029FF"
    "U+00002A00-U+00002AFF"
    "U+00002B00-U+00002BFF"
    "U+00002C00-U+00002CFF"
    "U+00002D00-U+00002DFF"
    "U+00002E00-U+00002EFF"
    "U+00002F00-U+00002FFF"
    "U+00003000-U+000030FF"
    "U+00003100-U+000031FF"
    "U+00003200-U+000032FF"
    "U+00003300-U+000033FF"
    "U+00003400-U+000034FF"
    "U+00003500-U+000035FF"
    "U+00003600-U+000036FF"
    "U+00003700-U+000037FF"
    "U+00003800-U+000038FF"
    "U+00003900-U+000039FF"
    "U+00003A00-U+00003AFF"
    "U+00003B00-U+00003BFF"
    "U+00003C00-U+00003CFF"
    "U+00003D00-U+00003DFF"
    "U+00003E00-U+00003EFF"
    "U+00003F00-U+00003FFF"
    "U+00004000-U+000040FF"
    "U+00004100-U+000041FF"
    "U+00004200-U+000042FF"
    "U+00004300-U+000043FF"
    "U+00004400-U+000044FF"
    "U+00004500-U+000045FF"
    "U+00004600-U+000046FF"
    "U+00004700-U+000047FF"
    "U+00004800-U+000048FF"
    "U+00004900-U+000049FF"
    "U+00004A00-U+00004AFF"
    "U+00004B00-U+00004BFF"
    "U+00004C00-U+00004CFF"
    "U+00004D00-U+00004DFF"
    "U+00004E00-U+00004EFF"
    "U+00004F00-U+00004FFF"
    "U+00005000-U+000050FF"
    "U+00005100-U+000051FF"
    "U+00005200-U+000052FF"
    "U+00005300-U+000053FF"
    "U+00005400-U+000054FF"
    "U+00005500-U+000055FF"
    "U+00005600-U+000056FF"
    "U+00005700-U+000057FF"
    "U+00005800-U+000058FF"
    "U+00005900-U+000059FF"
    "U+00005A00-U+00005AFF"
    "U+00005B00-U+00005BFF"
    "U+00005C00-U+00005CFF"
    "U+00005D00-U+00005DFF"
    "U+00005E00-U+00005EFF"
    "U+00005F00-U+00005FFF"
    "U+00006000-U+000060FF"
    "U+00006100-U+000061FF"
    "U+00006200-U+000062FF"
    "U+00006300-U+000063FF"
    "U+00006400-U+000064FF"
    "U+00006500-U+000065FF"
    "U+00006600-U+000066FF"
    "U+00006700-U+000067FF"
    "U+00006800-U+000068FF"
    "U+00006900-U+000069FF"
    "U+00006A00-U+00006AFF"
    "U+00006B00-U+00006BFF"
    "U+00006C00-U+00006CFF"
    "U+00006D00-U+00006DFF"
    "U+00006E00-U+00006EFF"
    "U+00006F00-U+00006FFF"
    "U+00007000-U+000070FF"
    "U+00007100-U+000071FF"
    "U+00007200-U+000072FF"
    "U+00007300-U+000073FF"
    "U+00007400-U+000074FF"
    "U+00007500-U+000075FF"
    "U+00007600-U+000076FF"
    "U+00007700-U+000077FF"
    "U+00007800-U+000078FF"
    "U+00007900-U+000079FF"
    "U+00007A00-U+00007AFF"
    "U+00007B00-U+00007BFF"
    "U+00007C00-U+00007CFF"
    "U+00007D00-U+00007DFF"
    "U+00007E00-U+00007EFF"
    "U+00007F00-U+00007FFF"
    "U+00008000-U+000080FF"
    "U+00008100-U+000081FF"
    "U+00008200-U+000082FF"
    "U+00008300-U+000083FF"
    "U+00008400-U+000084FF"
    "U+00008500-U+000085FF"
    "U+00008600-U+000086FF"
    "U+00008700-U+000087FF"
    "U+00008800-U+000088FF"
    "U+00008900-U+000089FF"
    "U+00008A00-U+00008AFF"
    "U+00008B00-U+00008BFF"
    "U+00008C00-U+00008CFF"
    "U+00008D00-U+00008DFF"
    "U+00008E00-U+00008EFF"
    "U+00008F00-U+00008FFF"
    "U+00009000-U+000090FF"
    "U+00009100-U+000091FF"
    "U+00009200-U+000092FF"
    "U+00009300-U+000093FF"
    "U+00009400-U+000094FF"
    "U+00009500-U+000095FF"
    "U+00009600-U+000096FF"
    "U+00009700-U+000097FF"
    "U+00009800-U+000098FF"
    "U+00009900-U+000099FF"
    "U+00009A00-U+00009AFF"
    "U+00009B00-U+00009BFF"
    "U+00009C00-U+00009CFF"
    "U+00009D00-U+00009DFF"
    "U+00009E00-U+00009EFF"
    "U+00009F00-U+00009FFF"
    "U+0000A000-U+0000A0FF"
    "U+0000A100-U+0000A1FF"
    "U+0000A200-U+0000A2FF"
    "U+0000A300-U+0000A3FF"
    "U+0000A400-U+0000A4FF"
    "U+0000A500-U+0000A5FF"
    "U+0000A600-U+0000A6FF"
    "U+0000A700-U+0000A7FF"
    "U+0000A800-U+0000A8FF"
    "U+0000A900-U+0000A9FF"
    "U+0000AA00-U+0000AAFF"
    "U+0000AB00-U+0000ABFF"
    "U+0000AC00-U+00
```



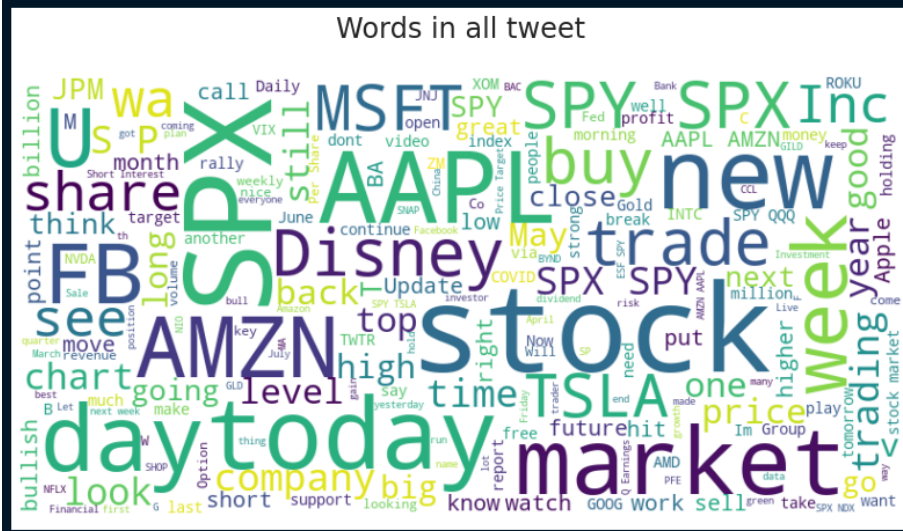
## Te pokrenuli te filtere

```
def arrange_text(ds):
    ds['text2'] = ds['text'].apply(emojify)
    ds['text2'] = ds['text2'].apply(handle)
    ds['text2'] = ds['text2'].apply(specialcode)
    ds['text2'] = ds['text2'].apply(hashtag)
    ds['text2'] = ds['text2'].apply(url)
    ds['text2'] = ds['text2'].apply(pic)
    ds['text2'] = ds['text2'].apply(html_tag)
    ds['text2'] = ds['text2'].apply(onlychar)
    ds['text2'] = ds['text2'].apply(decontracted)
    ds['text2'] = ds['text2'].apply(onlychar)
    ds['text2'] = ds['text2'].apply(tokenize_stem)
    ds['text2'] = ds['text2'].apply(remove_stopwords)
```

Prikazali smo sve riječi u tweet-ovima

```
> MI
words = ' '.join([text for text in data['text2']])
wordcloud = WordCloud(
    width=800, height=400, background_color='white', max_font_size=110)\
    .generate(words)

plt.figure(figsize=(14, 7))
plt.imshow(wordcloud, interpolation="bilinear")
plt.title('Words in all tweet\n', fontsize=24)
plt.axis('off')
plt.show()
```



Dalje, napisali smo funkciju koja će formirati word cloud po sentimentima koje želimo:

```
def wordcloud_by_sentiment(sentiment):
    not_ticker = []

    for text in data[data['sentiment']==sentiment]['text2']:
        for word in text.split():
            if word.upper() not in ticker_dic:
                not_ticker.append(word.lower())

    words = ' '.join([word for word in not_ticker])
    wordcloud = WordCloud(
        width=800, height=400, background_color='white', max_font_size=110, max_words=100).\
        generate(words)
    plt.figure(figsize=(14, 7))
    plt.imshow(wordcloud, interpolation="bilinear")
    plt.title('Words in ' + sentiment.capitalize() + ' tweets\n', fontsize=32)
    plt.axis('off')
    plt.show()
```

Te svako od nas ponaosob je formirao jedan od wordcloud-a za različite sentimente.

Ena:

```
wordcloud_by_sentiment('positive')
```

## Words in Positive tweets



Midhat:

```
wordcloud_by_sentiment('negative')
```

## Words in Negative tweets



Džemil:

```
wordcloud_by_sentiment('neutral')
```

### Words in Neutral tweets



Tabelarni prikaz algoritama sa naznakom ko je od članova tima zadužen za detaljnu analizu i implementaciju algoritma:

Student/metoda ML-a	Log. regression	Naive Bayes	SVM
Ena	<b>X</b>		
Džemil			<b>X</b>
Midhat		<b>X</b>	

Ena:

```

> *ML
#Logistic regression
from sklearn.linear_model import LogisticRegression

classifier = LogisticRegression()
classifier.fit(X_train, y_train)
score = classifier.score(X_test, y_test)

print("Accuracy:", score)
Accuracy: 0.7776

```

Džemil:

```

> *ML
#SVM
import numpy as np
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
clf_svm = make_pipeline(StandardScaler(with_mean=False), SVC(gamma='auto'))
clf_svm.fit(X_train, y_train)
clf_svm.score(X_test, y_test)
0.8144

```

Midhat:

```

#Naive Bayes
from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import CategoricalNB
from sklearn.naive_bayes import BernoulliNB

clf = BernoulliNB()
clf.fit(X_train, y_train)
score_clf = classifier.score(X_test, y_test)
print(score_clf)
0.7776

```

- 4) Detaljna analiza i Implementacija algoritama. Svaki član tima treba detaljno analizirati jedan algoritam-metodu i implementirati ga na osnovu tabele iz 3. Analiza znači sažeti opis metode sa teorijskog aspekta sa adekvatnom referencom na detaljni opis, analiza prihvatljivosti te metode za problem koji se rješava, implementacija metode – opis koji jezik ili okruženje se koristilo. Potrebno je priložiti kodove.

*Kodovi su predloženi u prošlom odjeljku. Korišten je SKLEARN paket jer u sebi sadrži sve metode koje smo odabrali da implementiramo. Implementirane metode izvršavaju svoje treniranje na sličan način. Glavna razlika je u tome koji modul je korišten. Razliku između metoda smo odabrali objasniti u teoretskoj formi zato jer priloženi kodovi su izuzetno slični.*

#### Ena:

Modelom logističke regresije opisujemo vezu između prediktora koji mogu biti neprekidni, binarni, kategorički, i kategoričke zavisne promjenljive. Na primjer zavisna promjenljiva može biti binarna - na osnovu nekih prediktora predviđamo da li će se nešto desiti ili ne. Ocjenjujemo zapravo vjerovatnoće pripadanja svakoj kategoriji za dat skup prediktora. U zavisnosti od tipa zavisne promjenljive imamo:

Binarnu logističku regresiju - zavisna promjenljiva je binarna. Na primjer: proći ili pasti test, odgovor da ili ne u anketi, visok ili nizak krvni pritisak

Nominalna logistička regresija - zavisna promjenljiva ima 3 ili više kategorija koje se ne mogu po vrednosti upoređivati. Na primjer: boje (crna, crvena, plava,...), internet pretraživač (Google, Bing, Yahoo!,...) ili u našem slučaju, pozitivan, neutralan i negativan sentiment.

Ordinalna logistička regresija - zavisna promjenljiva ima 3 ili više kategorija koje se prirodno mogu upoređivati, ali rangiranje ne mora da znači da su "rastojanja" između njih jednaka. Na primjer: ocjena (5-10), zdravstveno stanje (dobro, stabilno, ozbiljno, kritično) itd. Logistička regresija koristi kada zavisna promjenljiva uzima samo neki konačan skup vrijednosti. *Izvor: predavanja sa MU, RIMI*  
Ovo je slučaj kod nas i također, poznato je da ovaj oblik regresije je više robustan metod za klasifikaciju podataka kao što su tekst, što čini ovu metodu prihvatljivom u našem slučaju. Algoritam koji je implementiran je iz sklearn python paketa i kod je:

```
#Logistic regression
from sklearn.linear_model import LogisticRegression

classifier = LogisticRegression()
classifier.fit(X_train, y_train)
score = classifier.score(X_test, y_test)

print("Accuracy:", score)
Accuracy: 0.7776
```

Odakle se vidi da se inicijalizira objekat klase LogisticRegression koji je naš klasifikator. Pripremljene podatke (opisane u preprocessing dijelu) proslijeđujemo metodi fit ove klase. Na kraju ispitujemo tačnost naše metode preko test skupa pozivom metode score. Ispisujemo konačnu tačnost koja je jednaka 77.76%.

#### Midhat:

Naive Bayes je probabilistički klasifikator. Ovo znači da glavni mehanizam rada je usko povezan sa računanjem vjerovatnoće i uslovne vjerovatnoće događaja. Konkretno u našem slučaju, vjerovatnoće određene klase u odnosu na ulazne podatke. Naivni Bayes klasifikator se zasniva na Naive Bayes teoremi koja podrazumijeva gotovo potpuno neovisnost podatka. Pokazalo se da iako postoji određena ovisnost u podacima pri sentimentalnoj analizi, zaključujemo da ta ovisnost i korelacija nije dovoljno velika da utiče na tačnost NBC-a u tolikoj mjeri da metoda ne bude iskorištena. Odabrali smo metodu ne samo zbog toga,

već i zbog velike brzine treniranja i inferiranja testnog skupa. Generalno metode koje smo koristili su mnogo brže od njihovih pandana iz opšte, AI klase algoritama čija brzina treniranja je neuporedivo manja od ovih metoda, što konkretno za predviđanje koje je trenutno potrebno klijentu, se neće moći tolerisati za iole više korisnika i više podataka, što je konačni cilj ovakvog sistema u komercijalnom smislu.

Po Bayes-ovoj teoremi, uslovna vjerovatnoća se računa po sljedećoj formuli:

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})}$$

Te se klasifikacija vrši prema sljedećoj formuli:

$$\hat{y} = \operatorname{argmax}_{k \in \{1, \dots, K\}} p(C_k) \prod_{i=1}^n p(x_i | C_k)$$

Gdje je  $C_k$  klasa koju uzimamo i analiziramo. Ponaosob je potrebno analizirati i izračunati vjerovatnoću pripadnosti svakoj klasi i uzeti maksimum za rezultat. Pošto imamo samo tri klase, ovo je poprilično jednostavno uraditi i nisu potrebni veliki računarski resursi da bi se izvršila klasifikacija.

Potencijalni minus je što skup podataka postaje veći, to je veći proračun potreban kako bi se izvršila klasifikacija. Isto tako sa više podataka raste mogućnost da određeni podaci budu u jakoj korelaciji, koja smeta ovom klasifikatoru. Ovo su informacije koje je potrebno imati na umu pri implementaciji ovog klasifikatora. *Izvor: predavanja sa MU, RIMI*

Što se tiče implementacije, koristeći sklearn pakete za NaiveBayes smo testirali i odredili da BernoulliNB je najbolji klasifikator od tri ponuđena. Nakon inicijalizacije BernoulliNB klase bez parametara, prosljeđujemo trening podatke fit metodi ove klase. Nakon toga, koristeći score metodu zaključujemo da je tačnost 77.76%

```
#Naive Bayes
from sklearn.naive_bayes import GaussianNB
from sklearn.naive_bayes import CategoricalNB
from sklearn.naive_bayes import BernoulliNB

clf = BernoulliNB()
clf.fit(X_train, y_train)
score_clf = clf.score(X_test, y_test)
print(score_clf)

0.7776
```

Džemil:

SVM je namijenjen rješavanju binarnih klasifikacijskih problema ali također se može primjeniti i na multi-class klasifikacijski problem, što je nama potrebno u ovom slučaju. Konkretno, potrebna nam je troklasna klasifikacija za naše predstavljene klase (pozitivan, negativan i neutralan sentiment).

Poznato je da je SVM klasifikator nadogradnja MM koji zahtjeva potpuno razdvajanje klasa. Očigledno je da kod nas to nije slučaj jer određeni sentimentu su tranzientni kroz sve moguće klase (npr riječ stock). Očigledno je da je potrebno koristiti SVM klasifikator.

Pošto za SVM treniranje i testiranje ovisi isključivo od skalarnog proizvoda instanci koje se nalaze unutar i na margini, kao i skalarnog proizvoda instanci koje se nalaze na i unutar margine sa nepoznatim testnim skupom, ovaj metod je vrlo jednostavno implementirati sa obzirom na to koje metode su prethodno korištene (podaci su spremni, kao i rezultati za poređenje). *Izvor: predavanja sa MU, RIM1*

$$\max \sum_{i=1}^l \alpha_i - \frac{1}{2} \sum_{i=1}^l \sum_{j=1}^l \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$
$$f(\mathbf{x}^*) = \mathbf{w} \cdot \mathbf{u} + b = \sum_{i,j=1}^l \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}^* + b$$

Metod koji je korišten za implementaciju je iz paketa sklearn, konkretno SVC klasa svm modula unutar ovog paketa. Pošto je za SVM potrebno odrediti skaliranje podataka, napravljen je pipeline koristeći pipeline modul sklearn paketa u pythonu. Nakon formiranja pipeline-a, slično kao prethodne dvije metode, isti trenirng skup je dodijeljen fit metodi klasifikatora SVM-a sa default podacima koji su se pokazali dovoljnim, s obzirom na to da SVM klasifikator posjeduje optimizatore hiperparametara u tom slučaju koji će odabrati između ostalog i optimalnu funkciju kernela.

```
#SVM
import numpy as np
from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
clf_svm = make_pipeline(StandardScaler(with_mean=False), SVC(gamma='auto'))
clf_svm.fit(X_train, y_train)
clf_svm.score(X_test, y_test)

0.8144
```

Istrenirani model se evaluira preko testnog skupa i vidi se da je tačnost modela 81.44% što je i najveća tačnost dosad. Model izuzetno brzo dostavlja rezultat i brzo trenira, što ga čini mnogo više iskoristivim od možda RNN mreža za čije treniranje je potrebno mnogo vremena, a i vrijeme inferencije istih nije zavidno u odnosu na metode ML-a.

##### 5) Analiza dobijenih rezultata na podacima iz 2 upotrebom raznih metrika evaluacije. Analiza primjenjivosti sistema.

Rezultati koji su dobijeni su prihvatljivi s obzirom na to koje metode su iskorištene. Najvažnija metrika za evaluaciju je metrika tačnosti.

Metrika tačnosti ukazuje na to da se tačnost naših metoda kreće od 77-81%, što autori navode kao maksimalnu tačnost koju su i oni postigli.

Još jedna metrika koja je razmatrana je neovisnost klasa. Pri sentimentalnoj analizi je primijećeno da određene riječi su koliko-toliko isključivo vezane za određene klase.

Pri testiranju, ako su intenzivno korištene ove riječi, došlo je do pada tačnosti. Zaključujemo da ovisnost klasa o određenim riječima je umjereno visoka i to je nešto što se treba razmotriti ti povećavanju skupa podataka za treniranje ili testiranje.

Metrike za True Positive i True Negative su jedine koje imaju smisla biti iskorištene u ovom slučaju. False Negative i False Positive imaju smisla za dvoklasnu klasifikaciju, tako da se u ovom slučaju odbacuje.