



# **Sentimentalna analiza i predviđanje kretanja cijena dionica koristeći Twitter podatke**

Džemil Džigal  
Midhat Hodo  
Ena Pirija

# Dionice, indeksi, Calls & Puts



Predviđanje cijena pojedinih dionica je izuzetno resursno intenzivno i neisplativo

Predviđanje cijena indeksa je lakše jer isti objedinjuju više dionica u istu grupu

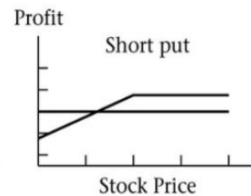
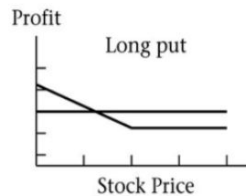
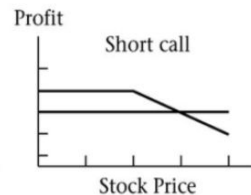
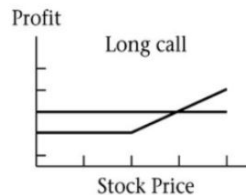
Calls - “klađenje” da će cijena dionice porasti

Put - “klađenje” da će cijena dionice pasti

Trajanje Call-ova i Put-ova. Kad zakazati strike-out? Long vs Short-term strike.

Long - “siguran ali mal profit”

Short - “rizičan ali visok profit”





# Uticaj socijalnih mreža na cijene dionica

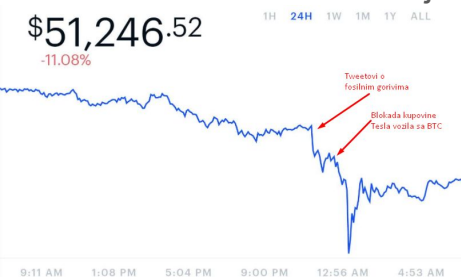
Na cijene dionica utiču ne samo standardni faktori kao stanje u industriji, potražnja i slično.

Utjecajni ljudi sa uspostavljenim dosegom na socijalnim mrežama utiču na cijene.

Elon Musk

Bitcoin Crash

Pad berze 2008 je plod rada ljudi, ne kolapsa industrije (hedge funds).



**Pinned Tweet**

**Elon Musk** @elonmusk

**You can now buy a Tesla with Bitcoin**

7:02 AM · Mar 24, 2021 · Twitter for iPhone

**Elon Musk** @elonmusk · 12s

To be clear, I strongly believe in crypto, but it can't drive a massive increase in fossil fuel use, especially coal

6 4 12

**Elon Musk** @elonmusk · 12m

It is high time there was a carbon tax!

16.1K 6.3K 38.5K



# Sentimentalna analiza Tweet podataka

Primijećen je uticaj tweet-ova na cijene dionica na berzi.

Ovi podaci se mogu analizirati kako bi se odredio sentiment kretanja cijene dionica.

Određene riječi u tweet-ovima su često bile povezane sa rastom ili padom cijena dionica.

Crash, bearish, sell, stop, short su često povezane sa negativnim sentimentom i padom dionica.

Boom, bullish, buy, start, good su često povezane sa pozitivnim sentimentom i rastom cijena.

Earning, trading, market, new su često povezane sa neutralnim sentimentom i održavanjem cijena.



# Sentimentalna analiza

Koristeći labele i vrijednosti u dataset-u, možemo vizualizirati sentiment i analizirati isti

```
1 def wordcloud_by_sentiment(sentiment):
2     not_ticker = []
3
4     for text in data[data['sentiment']==sentiment]['text2']:
5         for word in text.split():
6             if word.upper() not in ticker_dic:
7                 not_ticker.append(word.lower())
8
9     words = ' '.join([word for word in not_ticker])
10    wordcloud = WordCloud(
11        width=800, height=400, background_color='white', max_font_size=110, max_words=100).\
12        generate(words)
13    plt.figure(figsize=(14, 7))
14    plt.imshow(wordcloud, interpolation="bilinear")
15    plt.title('Words in ' + sentiment.capitalize() + ' tweets\n', fontsize=32)
16    plt.axis('off')
17    plt.show()
```





## Riječi sa negativnim sentimentom







# Priprema podataka

Klasična priprema tekstualnih podataka za NLP.

Uklanjanje UNICODE i specijalnih karaktera.

Koristiti emoji-e za sentiment ako sentiment nedostaje u podacima (NaN).

```
def arrange_text(ds):  
    ds['text2'] = ds['text'].apply(emoji)  
    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)
```



# Metode predikcije sentimenta

Logistička regresija - Jednostavan metod, brz i adekvatan za podatke koje koristimo.

Naive Bayes - Brz metod, jednostavan ali potencijalni problem uvode novi podaci (uvode ovisnost).

SVM - Robustan i brz metod, više otporan na probleme prethodnih metoda. Veća tačnost.



# Logistička regresija

Sklean paket, LogisticRegression biblioteka

```
1 #logistic regression
2 from sklearn.linear_model import LogisticRegression
3
4 classifier = LogisticRegression()
5 classifier.fit(X_train, y_train)
6 score = classifier.score(X_test, y_test)
7
8
9 print("Accuracy:", score)
```

Tačnost, 77.76%.



# Naive Bayes

Sklean paket, naive\_bayes modul, BernoulliNB klasa

```
1 #Naive Bayes
2 from sklearn.naive_bayes import GaussianNB
3 from sklearn.naive_bayes import CategoricalNB
4 from sklearn.naive_bayes import BernoulliNB
5
6 clf = BernoulliNB()
7 clf.fit(X_train, y_train)
8 score_clf = clf.score(X_test, y_test)
9 print(score_clf)
```

Tačnost, 77.76%.



# SVM

Sklean paket, svm modul, SVC klasa

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

Tačnost, 81.44%.



## Dalji rad i zaključak

RNN ? Da li vrijedi vremena treniranja za poboljšanje tačnosti od oko 5-10% (ako i toliko)?  
Vrijede li većeg vremena inferiranja?

Ako se implementira distribuirani sistem, performanse postaju vrlo važne jer se planira sugerisati kupovina **short-term** call-ova i put-ova na SNP500 indeks.

Umjesto RNN-ova, Fuzzy kontrolni sistemi? Manja tačnost, manje podataka ali veća brzina?

Koliko brzo je prebrzo?

Šta su prihvatljive performanse i rezultati?

Kod i dataset dostupni na:

[dzemildzical/DAP \(github.com\)](https://github.com/dzemildzical/DAP)



**Hvala na pažnji!**

**Q&A**