

# Analysis of political events using Twitter

Natural Language Processing Project#11

Marzie Faramarzzadeh

2210301

# Abstract

---

In this project, I retrieved 100MB of tweets related to Ukraine war from [Ukraine Conflict Twitter Dataset \(53.33M tweets\) | Kaggle](#). We performed sentiment analyzes on them using nltk package in python. Also, we continued processing these tweets using other NLP techniques such as: LDA clustering, WordCloud, NER, co-occurring words histogram, and empath categorization. Finding of the project revealed the potential of NLP in detecting how people feel (like/hate) about this war. We found that majorly people disliked this event and overall sentiment score is negative.

# Github page of the project

---

- <https://github.com/aiefaramarzzadeh/Natural-Language-Processing-project.git>



# Topic(s) investigated

---

- Sentiment analyzer of tweets..
- LDA (#topics =10 and #words =10) of operation.
- WordCloud representation
- Identifying the named-entities
- Co-occurring words with “Russia”
- Empath client categories
- Co-occurring words for “hate” and “like”
- Named-entity tagger and empath client categories for tweets containing modal verbs
- Named-entity tagger and empath client categories for tweets containing ‘wish’

# Relevant prior work

---

- MYKOLA et al. (2015) explored the use of the #SaveDonbassPeople hashtag (against the military operation in Eastern Ukraine turning Twitter into an online battleground). They found that Twitter was predominantly used as a propaganda outlet to broadcast opposing views on the ongoing conflict.
- Polyzos (2022) proposed the use of social media information as a real-time decision-making tool for significant events, using the war in Ukraine as a case study. He utilized the public's perception of the progression of events using sentiment analysis on 42 million tweets. He found that European currencies and markets experience an immediate negative response to conflict escalation “shocks”, while crude oil registers a delayed negative response. US stock markets seem unaffected, while the US Dollar responds positively to negative events of the war. His findings suggest that user generated content can be used as a decision-making tool when as important events unfold.
- After Russia invasion to Ukraine, a second battlefield has emerged in the online space, both in the use of social media to garner support for both sides of the conflict and also in the context of information warfare. In this paper, Chen et al. (2022 ) presented a collection of over 63 million tweets, from February 22, 2022 through March 8, 2022 that we are publishing for the wider research community to use (<https://github.com/echen102/ukraine-Russia>). Their preliminary analysis already showed evidence of public engagement with Russian state sponsored media and other domains that are known to push unreliable information.

# Data sources

- **Ukraine Conflict Twitter Dataset (54.26M tweets) by BwandoWando**
  - <https://www.kaggle.com/datasets/bwandoWando/ukraine-russian-crisis-twitter-dataset-1-2-m-rows>
- **Daily datasets of tweets about the ongoing Ukraine Russia Conflict**

## ACTIVITY STATS

VIEWS

**36887**

DOWNLOADS

**8154**

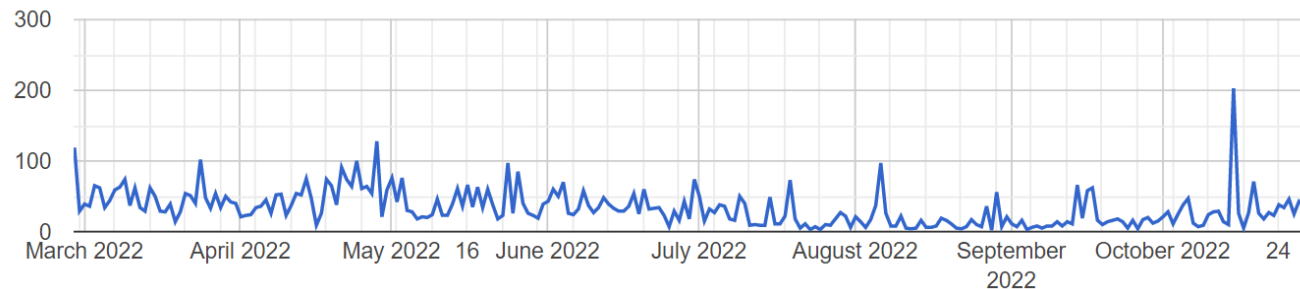
DOWNLOAD PER VIEW RATIO

**0.22**

TOTAL UNIQUE CONTRIBUTORS

**26**

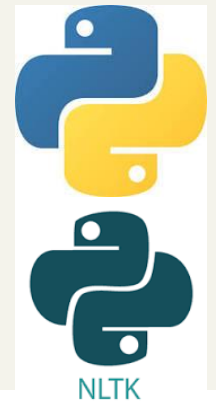
Downloads ▾



# Technologies and tools

---

- Python programming using Spyder IDE
- Packages used:
  - ✓ Spacy for NER
  - ✓ NLTK
  - ✓ Gensim and Pickle for LDA
  - ✓ Empath
  - ✓ Wordcloud



# Bag of words

## Raw tweet:

I support Ukraine ☺ 123...

I hate war #@! sdfa

.

.

.

## Clean tweet tokenized:

[I, support, ukraine]

[I, hate, war]

.

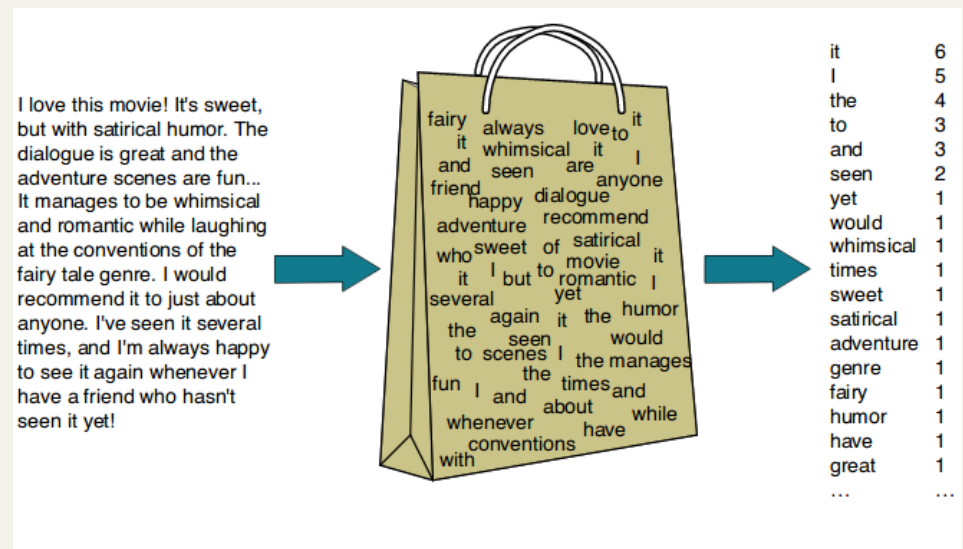
.

.

## Bag of word of clean tweets (excluding stopwords):

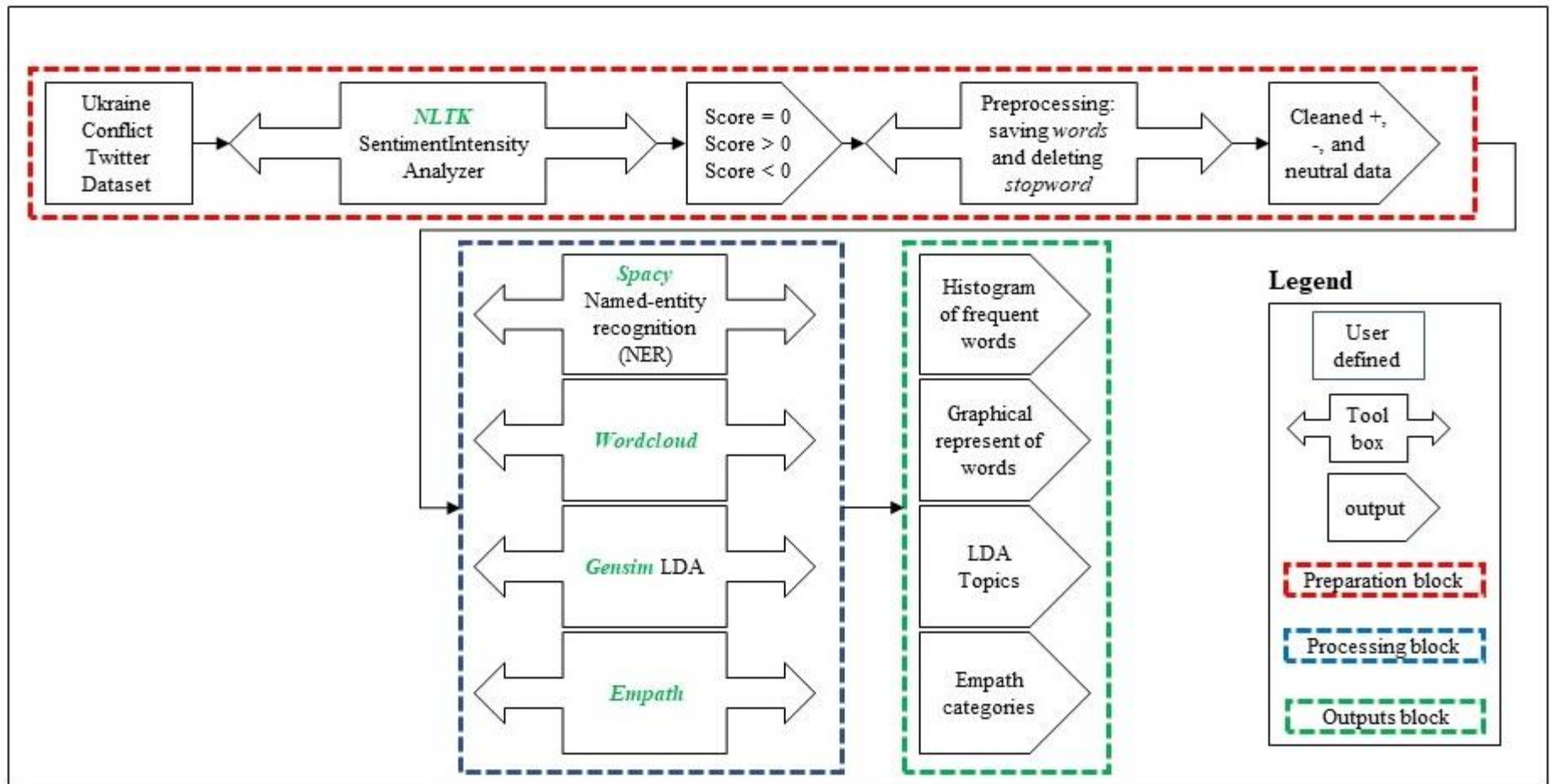
support ukraine hate war

...



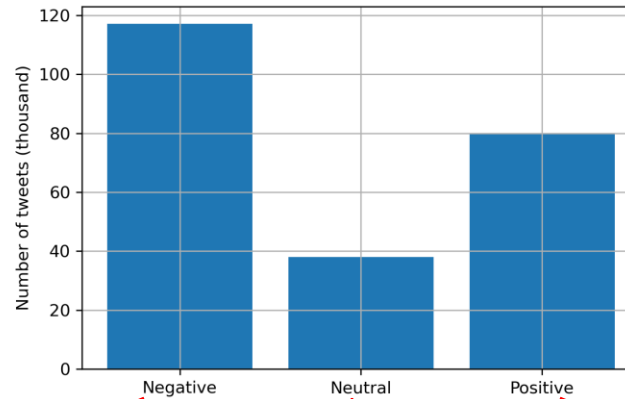


# Implementation details



For codes, see appendix

# Sentiment score



sent returned away son buried four former president  
exclusion zone news agency  
member public invasion via news  
state agency  
syndrome exclusion  
public council  
source member  
acute radiation  
council state  
radiation syndrome  
zone source  
agency exclusion  
seven suffering zone via  
old son told people moreover husband red forest room survive

politics try sam system tried make facility bombed respectively month  
spotted international two concentrate  
new collaboration target whether previously marked red  
perform new airport way jet drone  
international airport bam live  
concert perform make replace  
russia  
march concert  
first time  
bam bam closure border rare footage  
make world army bombed vehicle little leaving tea said loaded believe russia  
understand china

situation march long haul peace let giveaway  
companion know host follow tag attack win march find  
war stand  
one war  
peace nightingale  
mint day pick winner find government  
defence intelligence nightingale one  
mint peace  
stand mint  
tomorrow let winner mint



Age Group	Percentage
18-24	15%
25-34	20%
35-44	25%
45-54	20%
55-64	15%
65-74	10%
75-84	5%
85+	5%

[illegible][illegible][illegible][illegible]

politics try sam system large tried make facility board respectively cross  
spotted international international  
new collaboration international  
try make concentrate with time target whether previously marked  
replace understand airport way jet drop  
perform new bomb missile live  
city today said tried border ambassador cross new building location date  
international airport  
apparently target whether jet date unknown least days still  
concert perform make replace  
unknown location system use president push back phase two  
russia  
march concert  
together first missile unknown  
first time shah china  
world politics  
bam bam way march  
make world  
use apparently collaboration bam

[illegible]

# Negative dataset LDA topics

- ✓ 0.034\*"every" + 0.030\*"vehicle" + 0.030\*"drone" + 0.030\*"city" + 0.029\*"footage" + 0.026\*"shot" + 0.023\*"like" + 0.022\*"world" + 0.020\*"trucks" + 0.019\*"munitions"
- ✓ 0.037\*"former" + 0.032\*"home" + 0.029\*"returned" + 0.020\*"security" + 0.020\*"russia" + 0.018\*"part" + 0.013\*"according" + 0.013\*"seen" + 0.012\*"actually" + 0.011\*"underway"
- ✓ 0.068\*"war" + 0.048\*"russia" + 0.032\*"people" + 0.020\*"many" + 0.017\*"food" + 0.016\*"us" + 0.014\*"human" + 0.013\*"told" + 0.013\*"start" + 0.012\*"high"
- ✓ 0.024\*"un" + 0.021\*"russia" + 0.018\*"never" + 0.016\*"must" + 0.016\*"head" + 0.016\*"trump" + 0.013\*"war" + 0.013\*"northwest" + 0.011\*"recently" + 0.011\*"leave"
- ✓ 0.098\*"russia" + 0.029\*"war" + 0.020\*"use" + 0.017\*"us" + 0.012\*"military" + 0.011\*"tank" + 0.010\*"attack" + 0.010\*"west" + 0.009\*"personally" + 0.009\*"country"
- ✓ 0.034\*"region" + 0.032\*"sent" + 0.026\*"city" + 0.019\*"army" + 0.018\*"tank" + 0.017\*"already" + 0.015\*"berlin" + 0.014\*"note" + 0.014\*"offensive" + 0.013\*"see"
- ✓ 0.040\*"president" + 0.024\*"first" + 0.024\*"fighting" + 0.020\*"people" + 0.018\*"anonymous" + 0.017\*"vehicle" + 0.016\*"building" + 0.015\*"go" + 0.014\*"army" + 0.013\*"russia"
- ✓ 0.042\*"one" + 0.041\*"son" + 0.034\*"old" + 0.031\*"year" + 0.025\*"husband" + 0.023\*"got" + 0.023\*"moreover" + 0.021\*"house" + 0.017\*"least" + 0.017\*"several"
- ✓ 0.044\*"invasion" + 0.042\*"light" + 0.041\*"battle" + 0.040\*"shelling" + 0.038\*"get" + 0.029\*"fierce" + 0.028\*"traffic" + 0.026\*"say" + 0.025\*"city" + 0.024\*"worst"
- ✓ 0.107\*"zone" + 0.105\*"exclusion" + 0.104\*"agency" + 0.057\*"state" + 0.054\*"radiation" + 0.054\*"suffering" + 0.053\*"news" + 0.053\*"via" + 0.053\*"acute" + 0.052\*"syndrome"

# Neutral dataset LDA topics

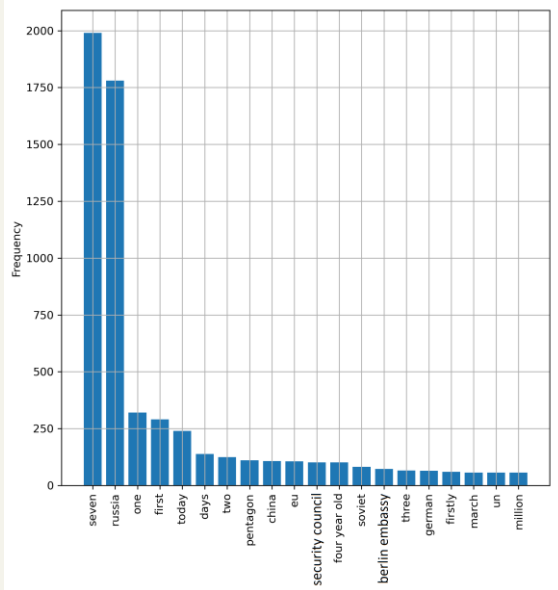
- ✓  $0.047 \times \text{"russia"} + 0.032 \times \text{"city"} + 0.031 \times \text{"today"} + 0.025 \times \text{"us"} + 0.019 \times \text{"east"} + 0.017 \times \text{"two"} + 0.016 \times \text{"still"} + 0.014 \times \text{"per"} + 0.013 \times \text{"bank"} + 0.012 \times \text{"mint"}$
- ✓  $0.058 \times \text{"unknown"} + 0.040 \times \text{"russia"} + 0.036 \times \text{"target"} + 0.033 \times \text{"use"} + 0.032 \times \text{"sam"} + 0.031 \times \text{"drone"} + 0.031 \times \text{"system"} + 0.029 \times \text{"location"} + 0.029 \times \text{"previously"} + 0.029 \times \text{"missile"}$
- ✓  $0.027 \times \text{"foreign"} + 0.027 \times \text{"minister"} + 0.026 \times \text{"people"} + 0.023 \times \text{"china"} + 0.018 \times \text{"million"} + 0.016 \times \text{"since"} + 0.016 \times \text{"working"} + 0.014 \times \text{"anonymous"} + 0.014 \times \text{"invasion"} + 0.013 \times \text{"real"}$
- ✓  $0.072 \times \text{"bam"} + 0.061 \times \text{"time"} + 0.052 \times \text{"new"} + 0.052 \times \text{"first"} + 0.049 \times \text{"march"} + 0.049 \times \text{"live"} + 0.043 \times \text{"international"} + 0.041 \times \text{"way"} + 0.040 \times \text{"spotted"} + 0.038 \times \text{"together"}$
- ✓  $0.052 \times \text{"army"} + 0.045 \times \text{"cross"} + 0.044 \times \text{"red"} + 0.039 \times \text{"building"} + 0.039 \times \text{"bombed"} + 0.029 \times \text{"south"} + 0.028 \times \text{"marked"} + 0.023 \times \text{"little"} + 0.019 \times \text{"price"} + 0.018 \times \text{"go"}$
- ✓  $0.109 \times \text{"russia"} + 0.036 \times \text{"day"} + 0.022 \times \text{"news"} + 0.020 \times \text{"video"} + 0.016 \times \text{"new"} + 0.015 \times \text{"tank"} + 0.014 \times \text{"war"} + 0.013 \times \text{"anti"} + 0.013 \times \text{"region"} + 0.013 \times \text{"china"}$
- ✓  $0.024 \times \text{"russia"} + 0.016 \times \text{"shot"} + 0.015 \times \text{"say"} + 0.015 \times \text{"end"} + 0.013 \times \text{"operation"} + 0.012 \times \text{"would"} + 0.012 \times \text{"house"} + 0.012 \times \text{"soldier"} + 0.011 \times \text{"water"} + 0.010 \times \text{"old"}$
- ✓  $0.027 \times \text{"russia"} + 0.026 \times \text{"days"} + 0.023 \times \text{"region"} + 0.022 \times \text{"via"} + 0.020 \times \text{"get"} + 0.020 \times \text{"part"} + 0.016 \times \text{"ago"} + 0.015 \times \text{"trump"} + 0.015 \times \text{"long"} + 0.013 \times \text{"armed"}$
- ✓  $0.085 \times \text{"russia"} + 0.035 \times \text{"breaking"} + 0.032 \times \text{"situation"} + 0.031 \times \text{"oil"} + 0.026 \times \text{"map"} + 0.018 \times \text{"approximate"} + 0.017 \times \text{"near"} + 0.014 \times \text{"war"} + 0.014 \times \text{"back"} + 0.013 \times \text{"gas"}$
- ✓  $0.089 \times \text{"make"} + 0.056 \times \text{"world"} + 0.047 \times \text{"china"} + 0.046 \times \text{"politics"} + 0.041 \times \text{"understand"} + 0.041 \times \text{"try"} + 0.041 \times \text{"tried"} + 0.040 \times \text{"replace"} + 0.039 \times \text{"border"} + 0.035 \times \text{"complete"}$

# Positive dataset LDA topics

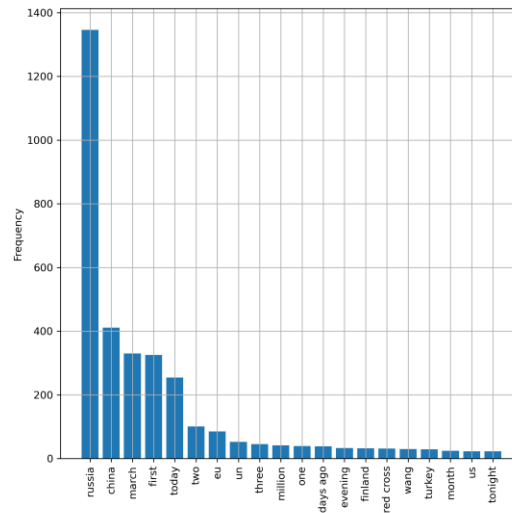
- ✓  $0.114 \cdot \text{"mint"} + 0.065 \cdot \text{"like"} + 0.062 \cdot \text{"know"} + 0.062 \cdot \text{"day"} + 0.050 \cdot \text{"tomorrow"} + 0.050 \cdot \text{"let"} + 0.048 \cdot \text{"giveaway"} + 0.045 \cdot \text{"tag"} + 0.045 \cdot \text{"excited"} + 0.044 \cdot \text{"winner"}$
- ✓  $0.029 \cdot \text{"help"} + 0.028 \cdot \text{"thank"} + 0.027 \cdot \text{"think"} + 0.018 \cdot \text{"russia"} + 0.017 \cdot \text{"time"} + 0.017 \cdot \text{"let"} + 0.015 \cdot \text{"good"} + 0.015 \cdot \text{"military"} + 0.015 \cdot \text{"even"} + 0.014 \cdot \text{"new"}$
- ✓  $0.058 \cdot \text{"russia"} + 0.026 \cdot \text{"win"} + 0.021 \cdot \text{"k"} + 0.021 \cdot \text{"form"} + 0.020 \cdot \text{"buy"} + 0.018 \cdot \text{"like"} + 0.016 \cdot \text{"people"} + 0.014 \cdot \text{"war"} + 0.014 \cdot \text{"us"} + 0.014 \cdot \text{"china"}$
- ✓  $0.043 \cdot \text{"intelligence"} + 0.042 \cdot \text{"government"} + 0.042 \cdot \text{"latest"} + 0.039 \cdot \text{"march"} + 0.038 \cdot \text{"situation"} + 0.037 \cdot \text{"defence"} + 0.036 \cdot \text{"response"} + 0.036 \cdot \text{"update"} + 0.035 \cdot \text{"find"} + 0.019 \cdot \text{"soon"}$
- ✓  $0.056 \cdot \text{"please"} + 0.029 \cdot \text{"support"} + 0.028 \cdot \text{"place"} + 0.028 \cdot \text{"world"} + 0.024 \cdot \text{"supporting"} + 0.023 \cdot \text{"thanks"} + 0.022 \cdot \text{"stop"} + 0.021 \cdot \text{"whole"} + 0.021 \cdot \text{"people"} + 0.020 \cdot \text{"army"}$
- ✓  $0.174 \cdot \text{"peace"} + 0.155 \cdot \text{"war"} + 0.152 \cdot \text{"one"} + 0.142 \cdot \text{"mint"} + 0.141 \cdot \text{"stand"} + 0.136 \cdot \text{"nightingale"} + 0.006 \cdot \text{"promise"} + 0.003 \cdot \text{"nice"} + 0.003 \cdot \text{"invasion"} + 0.003 \cdot \text{"difficult"}$
- ✓  $0.057 \cdot \text{"russia"} + 0.035 \cdot \text{"morning"} + 0.026 \cdot \text{"series"} + 0.019 \cdot \text{"peace"} + 0.017 \cdot \text{"used"} + 0.017 \cdot \text{"air"} + 0.016 \cdot \text{"system"} + 0.016 \cdot \text{"city"} + 0.016 \cdot \text{"military"} + 0.015 \cdot \text{"agreement"}$
- ✓  $0.049 \cdot \text{"tell"} + 0.033 \cdot \text{"neo"} + 0.027 \cdot \text{"truth"} + 0.026 \cdot \text{"show"} + 0.022 \cdot \text{"time"} + 0.022 \cdot \text{"german"} + 0.019 \cdot \text{"never"} + 0.019 \cdot \text{"back"} + 0.019 \cdot \text{"like"} + 0.019 \cdot \text{"russia"}$
- ✓  $0.036 \cdot \text{"army"} + 0.035 \cdot \text{"good"} + 0.024 \cdot \text{"get"} + 0.022 \cdot \text{"people"} + 0.016 \cdot \text{"military"} + 0.016 \cdot \text{"gas"} + 0.015 \cdot \text{"payment"} + 0.014 \cdot \text{"little"} + 0.014 \cdot \text{"make"} + 0.013 \cdot \text{"photo"}$
- ✓  $0.058 \cdot \text{"russia"} + 0.022 \cdot \text{"release"} + 0.022 \cdot \text{"project"} + 0.021 \cdot \text{"community"} + 0.020 \cdot \text{"giving"} + 0.020 \cdot \text{"follow"} + 0.018 \cdot \text{"tag"} + 0.018 \cdot \text{"help"} + 0.018 \cdot \text{"celebrate"} + 0.018 \cdot \text{"mon"}$

# NER for each dataset

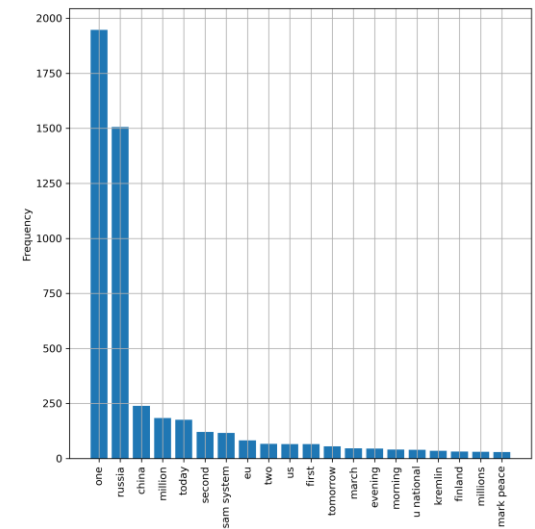
## Negative



## Neutral



## Positive

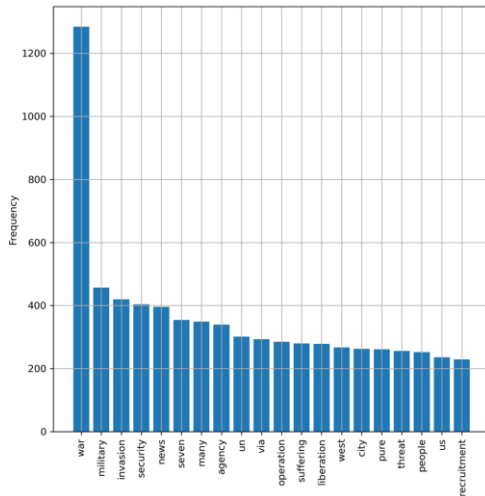


Named-entity can be a country, money, organization, months, cardinal numbers, ...

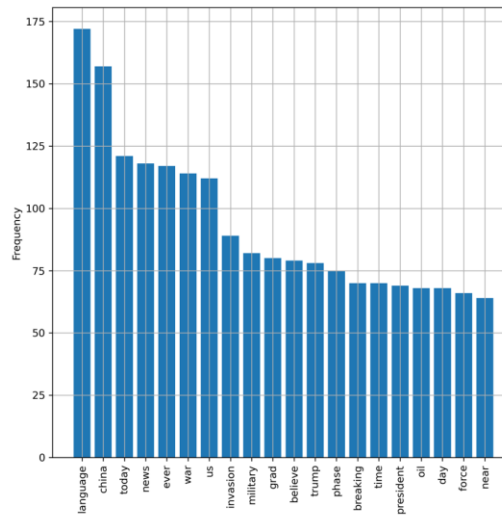


# Co-occurring words with Russia

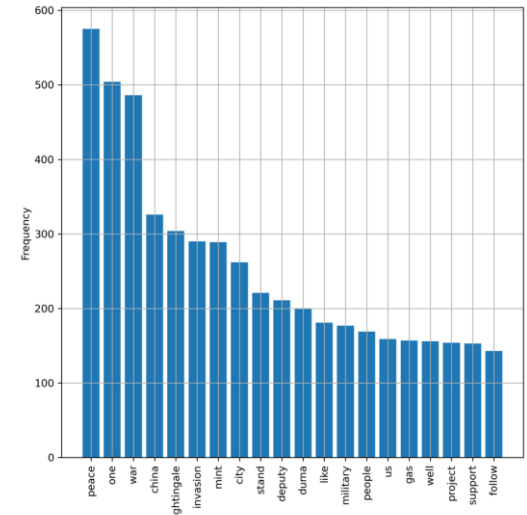
## Negative



## Neutral



## Positive





# Empath categories

---

## Negative:

war  
death  
government  
leader  
fight  
work  
banking  
violence  
military  
suffering  
pain  
dominant\_heirarchical  
politics  
law  
weapon  
hate  
sadness  
health  
terrorism  
anger  
poor  
royalty  
kill  
crime  
children

## Neutral:

war  
government  
military  
traveling  
musical  
leader  
music  
air\_travel  
vacation  
dominant\_heirarchical  
law  
politics  
speaking  
meeting  
communication  
economics  
party  
rural  
tourism  
business  
money  
home  
terrorism  
celebration  
weapon  
anonymity

## Positive:

war  
money  
fight  
smell  
death  
military  
body  
leader  
giving  
communication  
government  
party  
help  
celebration  
business  
achievement  
optimism  
positive\_emotion  
speaking  
dominant\_heirarchical  
competing  
play  
negative\_emotion  
politics  
order  
gain

# Calendar

---

- All *specifications* are done
- *Github* page is created, and codes uploaded in:
  - ✓ <https://github.com/aiefaramarzzadeh/Natural-Language-Processing-project.git>
- Writing *report*: 1-4 November 2022

# Questions

---

Thanks





# Appendix

# Sentiment score

```
def sentiment_score(self, tweets):
    from numpy import sign
    import nltk
    nltk.download('vader_lexicon')
    from nltk.sentiment import SentimentIntensityAnalyzer

    sia = SentimentIntensityAnalyzer() #calculating sentiment score

    scores = [sign(sia.polarity_scores(tweet)['compound']) for tweet in tweets]

    return scores

def sentiment_hist(self, scores):
    import numpy as np
    import matplotlib.pyplot as plt

    labels, counts = np.unique(scores, return_counts=True) #plotting histogram
    plt.bar(['Negative', 'Neutral', 'Positive'], (counts/1000), align='center')
    plt.grid()
    plt.ylabel('Number of tweets (thousand)')
    plt.savefig('histogram.png', dpi = 300)
```

# Cleaning + LDA

```
def pre_processor(self, text):
    from spacy.lang.en import English
    import nltk
    #nltk.download('stopwords')
    #nltk.download('words')

    en_stop = set(nltk.corpus.stopwords.words('english'))
    words = set(nltk.corpus.words.words())

    parser = English()

    tokens = list(parser(text))

    trash = ['#', '@', ';', ',', '.', ':', '/', '%', '$', '!', '*', '&', ')', '(', '_']
    trash.extend(list(en_stop))

    clear_tokenized_tweet = [str(elem).lower() for elem in tokens
                             if str(elem).lower() in words and str(elem).lower() not in trash]
```

pre-  
processing

```
import pandas as pd
import gensim
from gensim import corpora
import ast

tokenized_tweets = [ast.literal_eval(elem) for elem in tokenized_tweets]

dictionary = corpora.Dictionary(tokenized_tweets)
corpus = [dictionary.doc2bow(text) for text in tokenized_tweets]

import pickle
pickle.dump(corpus, open('corpus.pkl', 'wb'))
dictionary.save('dictionary.gensim')

ldamodel = gensim.models.ldamodel.LdaModel(corpus, num_topics = topics_num,
                                             id2word = dictionary, passes=15)
ldamodel.save('model10.gensim')

topics = ldamodel.print_topics(num_words = words_num)

df = pd.DataFrame(topics)
df.to_csv('task2_topics_'+name+'.csv')
```

LDA

# WordCloud

---

```
wordcloud_clear = WordCloud(width = 800, height = 800,  
                             background_color = 'white',  
                             min_font_size = 10).generate(words_clean)  
  
wordcloud_dirty = WordCloud(width = 800, height = 800,  
                             background_color = 'white',  
                             min_font_size = 10).generate(words_dirty)  
  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud_clear)  
plt.axis("off")  
plt.tight_layout(pad = 0)  
plt.savefig('wordcloud_clear_'+name+'.png', dpi = 300)  
  
plt.figure(figsize = (8, 8), facecolor = None)  
plt.imshow(wordcloud_dirty)  
plt.axis("off")  
plt.tight_layout(pad = 0)  
plt.savefig('wordcloud_dirty_'+name+'.png', dpi = 300)
```

# NER

---

```
named_entities = NER(words_clean)
ner_tweets = named_entities.ents
words = [elem.text for elem in ner_tweets]
labels = [elem.label_ for elem in ner_tweets]
spacy.explain("GPE") #explains what is each label of NER

words = pd.Series(nltk.FreqDist(words))
words = words.sort_values(ascending=False)

plt.figure(figsize=(8,8), dpi=300)
plt.bar(words.index[0:20], words.values[0:20])
plt.xticks(rotation=90)
plt.grid()
```



# N-gram

---

```
words_cnetered = [x for x in list(ngrams(words_clean.split(),7)) if x[3] == keyword]

words1 = [" ".join(elem) for elem in words_cnetered]
words = ' '.join([str(elem) for elem in words1])

words = pd.Series(nltk.FreqDist(words.split()))
words = words.sort_values(ascending=False)

plt.figure(figsize=(8,8), dpi=300)
plt.bar(words.index[1:21], words.values[1:21])
plt.xticks(rotation=90)
plt.grid()
plt.ylabel('Frequency')
plt.savefig('task5_'+keyword+'_freq_'+name+'.png', dpi = 300)
```

# Empath

```
def empath(self, input_clear, name):
    from empath import Empath
    import pandas as pd
    import ast
    lexicon = Empath()

    words_clean1 = [" ".join(ast.literal_eval(elem)) for elem in input_clear]
    words_clean = ' '.join([str(elem) for elem in words_clean1])

    categories = pd.Series(lexicon.analyze(words_clean, normalize=True))

    categories = categories[categories.values > 0]
    categories = categories.sort_values(ascending=False)

    categories.to_csv('task6_'+name+'_empath.csv')
```

['I','love','artificial','intelligence'] => ['I love artificial intelligence']  
['I','support','ukraine'] => ['I support ukraine']

love artificial  
intelligence.  
support ukraine

# Finding synonym

---

```
from nltk.corpus import wordnet

synonyms = []

for syn in wordnet.synsets(word):
    for syn in syn.lemmas():
        synonyms.append(syn.name())

synonyms = set(synonyms)
return synonyms
```

This can be like/hate

