Reyhane Shahrokhian
99521361
r_shahrokhian@comp.iust.ac.ir

Project Phaze1 Report
Natural Language Processing

2023-06-06

# 1 Gathering data

Due to the fact that tweets older than a week cannot be crawled using the official Twitter APIs, the source of data is actually the Telegram "Farsi Twitter" channel. So, crawling data is a bit different and it can't be done by a script or something like that.In fact i use the Export-Chat feature of telegram. After doing that , you will get json files that are all in data/row folder of the project.

# 2 Preprocessing and Labeling

To preprocess the data(cleaning data, word braking, sentence breaking, sentiment analysing) you should run related cells of project_phase1.ipynb which is in src folder. After running them you will get some files in your drive that are all in data/clean, data/wordbroken and data/sentencebroken folders of the project.
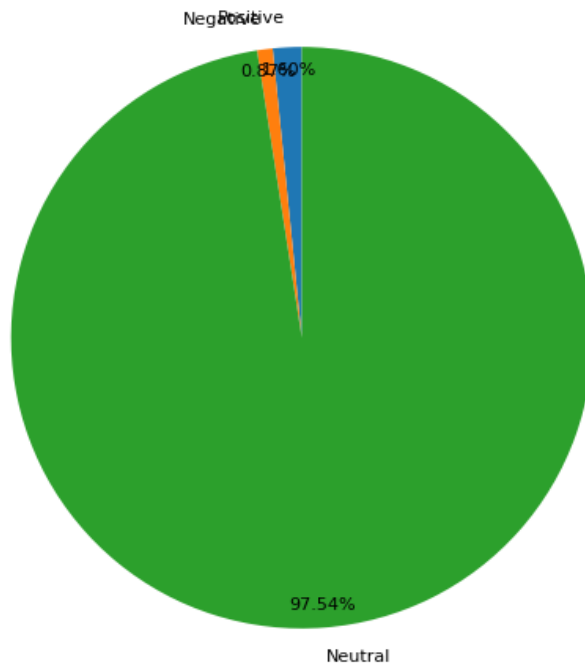
# 3 Statics

To show the static details(Data count of each label, word and sentence counts of labels, unique word count of labels, common and uncommon counts of words in comparison of labels, top10 words of each label and all words, top10 words of labels according to RNF and TF-IDF) you could run related cells of project_phase1.ipynb(statics section) which is in src folder. After running them you will get tables as csv file and charts as png in your drive that are all in stats folders of the project.
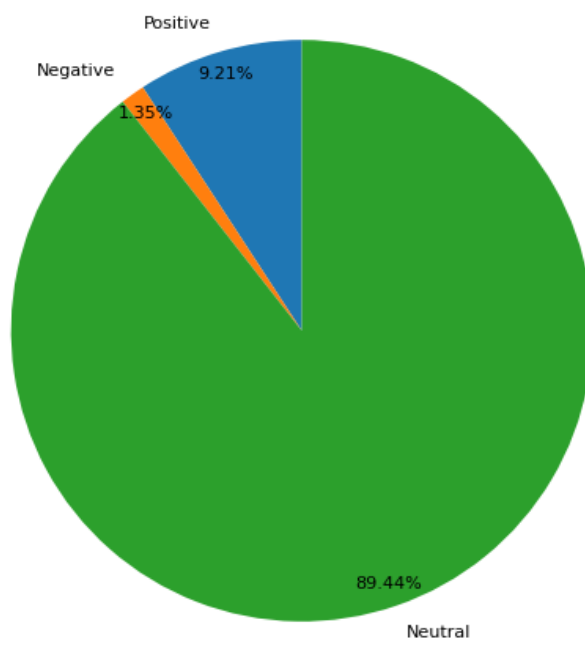
## 3.1 Statics of current data:

| Label | Data count |
|---|---|
| 2017Negative | 39 |
| 2017Positive | 72 |
| 2017Neutral | 4396 |
| 2020Negative | 69 |
| 2020Positive | 469 |
| 2020Neutral | 4555 |
| 2023Negative | 67 |
| 2023Positive | 550 |
| 2023Neutral | 4995 |

| Label | word count | unique word count | sentence count |
|---|---|---|---|
| 2017 | 79058 | 17212 | 5203 |
| 2020 | 118308 | 20444 | 6274 |
| 2023 | 194177 | 28396 | 9842 |

## sentiment analyse of 2017

Negative Positive
0.87% 1.60%

97.54%

Neutral

## sentiment analyse of 2020

Positive
Negative 9.21%
1.35%

89.44%

Neutral

## sentiment analyse of 2023

Positive

Negative

9.80%

1.19%

89.01%

Neutral

## 3.2 Top10 words of each label:

### 3.2.1 2017:

| Word | Count |
|------|-------|
|      | 259   |
|      | 254   |
|      | 251   |
|      | 250   |
|      | 242   |
|      | 240   |
|      | 233   |
|      | 221   |
| "    | 221   |
|      | 220   |
|      | 215   |

Line Chart of unique words 2017

### 3.2.2 2020:

| Word | Count |
|------|-------|
|      | 1151  |
|      | 1132  |
|      | 1083  |
|      | 1079  |
|      | 836   |
|      | 677   |
|      | 672   |
| !    | 653   |
|      | 493   |
|      | 436   |



Line Chart of unique words 2020

### 3.2.3 2023:

| Word | Count |
|------|-------|
|      | 927   |
|      | 907   |
|      | 864   |
|      | 687   |
|      | 660   |
|      | 649   |
|      | 616   |
|      | 586   |
|      | 582   |
|      | 573   |
|      | 568   |



Line Chart of unique words 2023

## 3.3 Common and Uncommon counts of words:

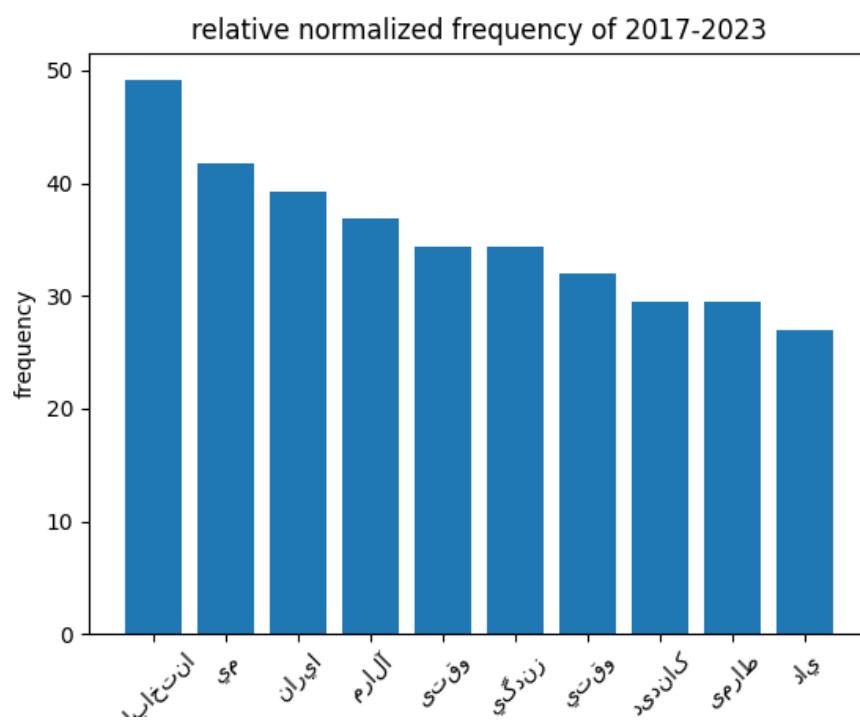| Label     | common word | uncommon word |
|-----------|-------------|---------------|
| 2017-2020 | 7439        | 9773          |
| 2017-2023 | 8500        | 8712          |
| 2020-2017 | 7439        | 13005         |
| 2020-2023 | 9522        | 10922         |
| 2023-2017 | 8500        | 19896         |
| 2023-2020 | 9522        | 18874         |

## Common Word Counts



## Uncommon Word Counts

## 3.4 Relative Normalized Frequency(RNF):

### 3.4.1 2017-2020 Common words by RNF

relative normalized frequency of 2017-2020



### 3.4.2 2017-2023 Common words by RNF

relative normalized frequency of 2017-2023

### 3.4.3 2020-2017 Common words by RNF



relative normalized frequency of 2020-2017

### 3.4.4 2020-2023 Common words by RNF



relative normalized frequency of 2020-2023

### 3.4.5    2023-2017 Common words by RNF
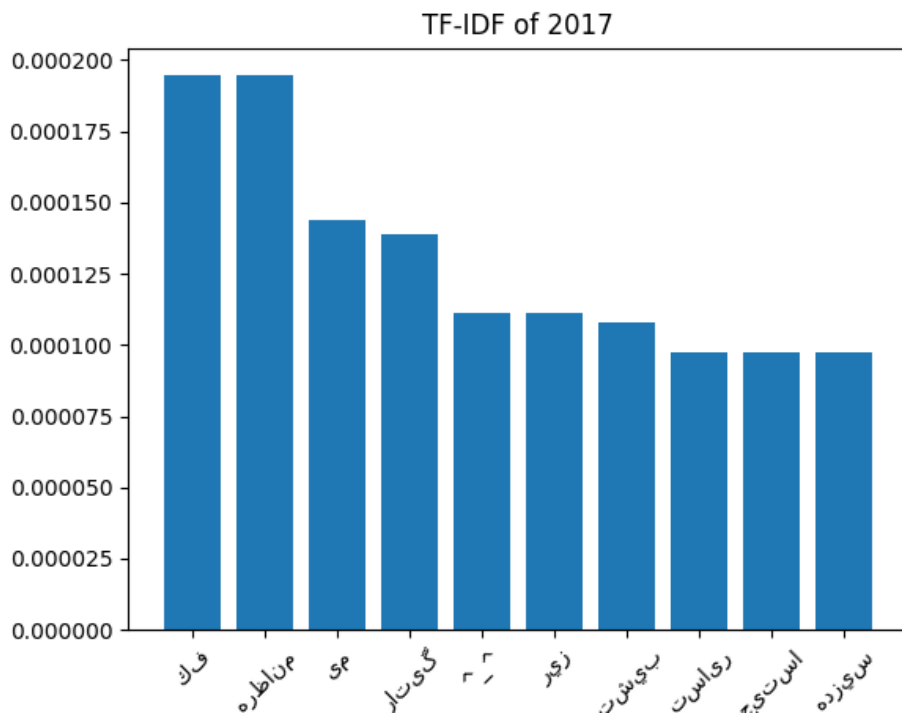


### 3.4.6    2023-2020 Common words by RNF
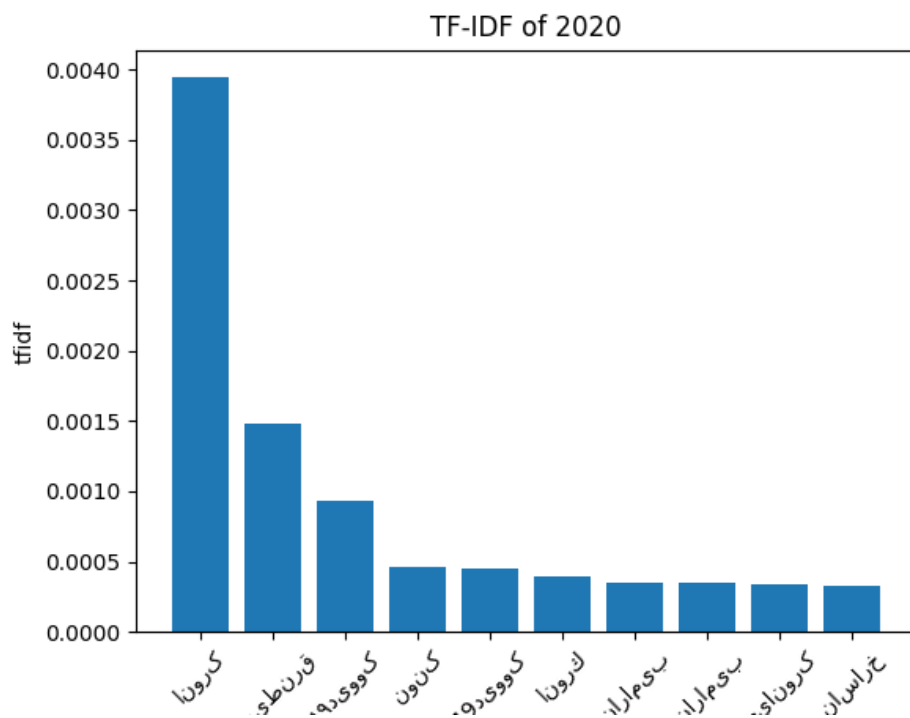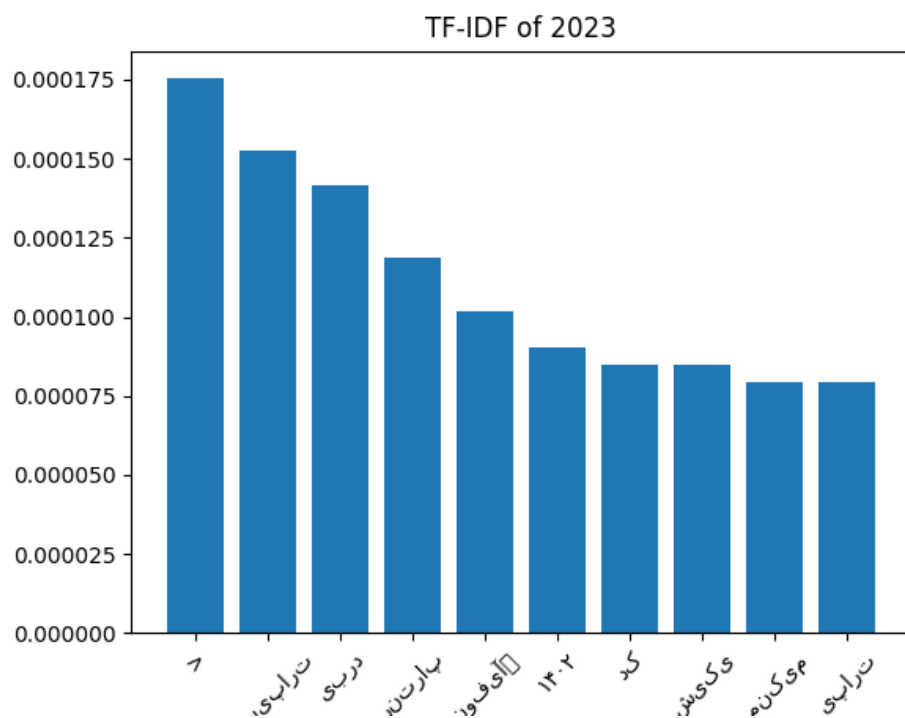
## 3.5 TF-IDF:

### 3.5.1 2017 tf-idf most Common words



### 3.5.2 2020 tf-idf Common words

### 3.5.3 2023 tf-idf most Common words



## 3.6 histogram of all word:

| Word | Count |
|---|---|
| | 1446 |
| | 1290 |
| | 1274 |
| | 1211 |
| | 1173 |
| | 1140 |
| | 1126 |
| | 1071 |
| | 1001 |
| | 963 |
| | 959 |

Line Chart of all words