



چالش دوم: Tweet Sentiment Extraction

فراہم آورنده: سینا حیدری

زمستان ۹۸

برای مشاهده‌ی جزئیات این چالش در کگل به لینک زیر مراجعه کنید.

<https://www.kaggle.com/c/tweet-sentiment-extraction/overview>

۱ در مورد این چالش

"My ridiculous dog is amazing." [sentiment: positive]

به دلیل تعداد بالای توییت‌هایی که در هر لحظه ردوبدل می‌شود، تایین اینکه تمایل یک توییت مشخص بر روی شرکت‌ها، اشخاص یا برنده‌ها چه تاثیری (مثبت یا منفی) می‌گذارد سخت است. امروزه، پی بردن به تمایل متن‌ها از آن جهت اهمیت یافته که تصمیم‌ها و واکنش‌ها در طی چند لحظه اتفاق افتاده و به‌روزرسانی می‌شود. اما چه کلماتی باعث تعیین تمایل متن‌ها می‌شود؟ در این چالش، باید بخشی از یک توییت (کلمه یا عبارت) را جدا کنیم که تمایل را تحت تاثیر قرار می‌دهد.

۲ راه حل

در مورد این گزارش

در این نوت‌بوک، ساختار مجموعه داده را به طور خلاصه توضیح داده، ویژگی‌های جدید را تولید و آنالیز می‌کنیم؛ سپس، مجموعه داده را با استفاده از Matplotlib، Seaborn و Plotly مصورسازی خواهیم کرد. ما به این مسئله، به عنوان یک مسئله‌ی NER^۱ خواهیم پرداخت. بررسی اکتشافی داده، اطلاعاتی به ما می‌دهد که بر اساس آن می‌توانیم نوع مدل خود را انتخاب کرده و به‌درستی شخصی سازی کنیم.

Named Entity Recognition

NER شاید اولین قدم برای استخراج اطلاعات باشد که تلاش می‌کند موجودیت‌های اسم‌دار را در متن شناسایی کرده و آن‌ها را به دسته‌های از پیش تعیین شده (همچون اسم افراد، سازمان‌ها، موقعیت‌ها، عبارت‌های زمان، کمیت‌ها، واحدهای پول، درصدها و غیره) نسبت دهد. NER در زمینه‌های مختلف NLP^۲ کاربرد دارد و می‌تواند به‌سوال‌های دنیای واقعی جواب دهد. نمونه‌هایی از این سوال‌ها را در زیر می‌بینید:

^۱ Named-entity recognition
^۲ Natural Language Processing

○ کدام شرکت‌ها در مقاله‌ی خبری مشخصی ذکر شده بودند؟

○ آیا محصول مشخصی در شکایت‌ها یا نقدها ظاهر شده است؟

○ آیا یک تویییت شامل اسم یک شخص می‌شود؟ آیا این تویییت موقعیت شخص را هم شامل می‌شود؟

ما از Spacy برای ساختن مدل‌های شخصی‌سازی شده‌مان استفاده خواهیم برد. نمونه‌یی از کاربرد NER بر اساس مدل اصلی SpaCy که با مجموعه‌داده‌ی OntoNotes 5 آموزش دیده را در زیر مشاهده می‌کنید. این مدل از موجودیت‌های زیر پشتیبانی می‌کند.

| TYPE | DESCRIPTION |
|-------------|------------------------------------------------------|
| PERSON | People, including fictional. |
| NORP | Nationalities or religious or political groups. |
| FAC | Buildings, airports, highways, bridges, etc. |
| ORG | Companies, agencies, institutions, etc. |
| GPE | Countries, cities, states. |
| LOC | Non-GPE locations, mountain ranges, bodies of water. |
| PRODUCT | Objects, vehicles, foods, etc. (Not services.) |
| EVENT | Named hurricanes, battles, wars, sports events, etc. |
| WORK_OF_ART | Titles of books, songs, etc. |
| LAW | Named documents made into laws. |
| LANGUAGE | Any named language. |
| DATE | Absolute or relative dates or periods. |
| TIME | Times smaller than a day. |
| PERCENT | Percentage, including "%". |
| MONEY | Monetary values, including unit. |
| QUANTITY | Measurements, as of weight or distance. |
| ORDINAL | "first", "second", etc. |
| CARDINAL | Numerals that do not fall under another type. |

طبقه‌بندی در سطح موجودیت، روی متن نمونه:

```

1 doc = nlp('European authorities fined Google a record $5.1 billion on
  Wednesday for abusing its power in the mobile phone market and
  ordered the company to alter its practices')
2 pprint([(X.text, X.label_) for X in doc.ents])
3

```

```

[('European', 'NORP'),
 ('Google', 'ORG'),
 ('$5.1 billion', 'MONEY'),
 ('Wednesday', 'DATE')]

```

در بخش بررسی اکتشافی داده خواهیم دید که برای بیشتر توییت‌های بی‌تمایل (neutral) جواب (selected_text) دقیقاً برابر با متن توییت می‌باشد؛ از این رو مدلی برای این دسته از توییت‌ها تعلیم نخواهیم داد. اما دو دسته توییت‌های منفی و مثبت وجود دارد که هر کدام نیازمند یک مدل هستند. به‌طور خلاصه، دو موجودیت اسم‌دار (مثبت و منفی) خواهیم داشت که برای هر کدام یک مدل آموزش می‌دهیم. ورودی مدل توییت‌های مثبت، آن سطرهایی از مجموعه‌داده است که تمایل مثبت دارد و ورودی مدل توییت‌های منفی آن سطرهایی از مجموعه‌داده است که تمایل منفی دارند. عبارت یا کلماتی که توسط NER شناسایی می‌شود نیز همان selected_text می‌باشد.

۳ پیاده‌سازی

بارگذاری پکیج‌های مورد نیاز

```
1 import re
2 import string
3 import numpy as np
4 import random
5 import pandas as pd
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 %matplotlib inline
9 from plotly import graph_objs as go
10 import plotly.express as px
11 import plotly.figure_factory as ff
12 from collections import Counter
13
14 from PIL import Image
15 from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
16
17
18 import nltk
19 from nltk.corpus import stopwords
20
21 from tqdm import tqdm
22 import os
23 import nltk
24 import spacy
25 import random
26 from spacy.util import compounding
27 from spacy.util import minibatch
28
29 import warnings
30 warnings.filterwarnings("ignore")
31
32 import os
33 for dirname, _, filenames in os.walk('/kaggle/input'):
34     for filename in filenames:
35         print(os.path.join(dirname, filename))
36
37 # Any results you write to the current directory are saved as output.
```

خواندن داده

```
1 print(train.shape)
2 print(test.shape)
3
```

ما ۲۷۴۸۶ توییت در مجموعه‌ی آموزشی و ۳۵۳۵ توییت در داده‌ی تست داریم.

```
(27481, 4)
(3534, 3)
```

```
1 train.info()
2
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 27481 entries, 0 to 27480
Data columns (total 4 columns):
textID      27481 non-null object
text        27480 non-null object
selected_text 27480 non-null object
sentiment    27481 non-null object
dtypes: object(4)
memory usage: 858.9+ KB
```

ما تنها یک مقدار null در داده‌ی آموزشی داریم. سطر شامل این مقدار را از داده‌ی آموزشی حذف می‌کنیم.

```
1 train.dropna(inplace=True)
2
```

```
1 train.info()
2
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 3534 entries, 0 to 3533
Data columns (total 3 columns):
textID      3534 non-null object
text        3534 non-null object
sentiment    3534 non-null object
dtypes: object(3)
memory usage: 83.0+ KB
```

هیچ مقدار null در داده‌ی تست وجود ندارد.

بررسی اکتشافی داده

```
1 train.head()
2
```

Out[8]:

| | textID | text | selected_text | sentiment |
|---|------------|-------------------------------------------------|-------------------------------------|-----------|
| 0 | cb774db0d1 | I'd have responded, if I were going | I'd have responded, if I were going | neutral |
| 1 | 549e992a42 | Sooo SAD I will miss you here in San Diego!!! | Sooo SAD | negative |
| 2 | 088c60f138 | my boss is bullying me... | bullying me | negative |
| 3 | 9642c003ef | what interview! leave me alone | leave me alone | negative |
| 4 | 358bd9e861 | Sons of ***, why couldn't they put them on t... | Sons of ***, | negative |

```
1 train.describe()
2
```

Out[9]:

| | textID | text | selected_text | sentiment |
|--------|------------|--------------------------------------------------|---------------|-----------|
| count | 27480 | 27480 | 27480 | 27480 |
| unique | 27480 | 27480 | 22463 | 3 |
| top | 609f4a0832 | can't wait to crack it open and no doubt will... | good | neutral |
| freq | 1 | 1 | 199 | 11117 |

حال، بیاید نگاهی به توزیع توییت‌ها در داده‌ی تست بی‌اندازیم.

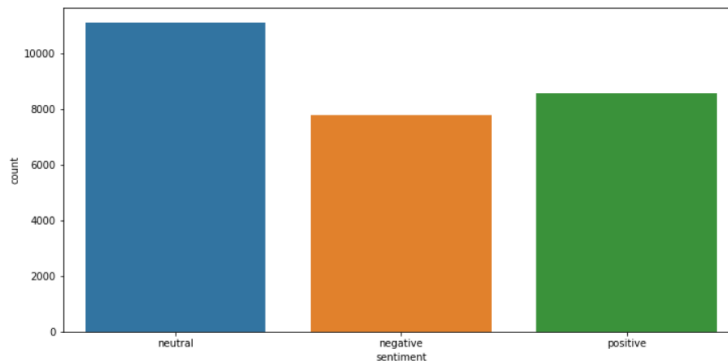
```
1 temp = train.groupby('sentiment').count()['text'].reset_index().
2 sort_values(by='text',ascending=False)
3 temp.style.background_gradient(cmap='Purples')
```

Out[10]:

| | sentiment | text |
|---|-----------|-------|
| 1 | neutral | 11117 |
| 2 | positive | 8582 |
| 0 | negative | 7781 |

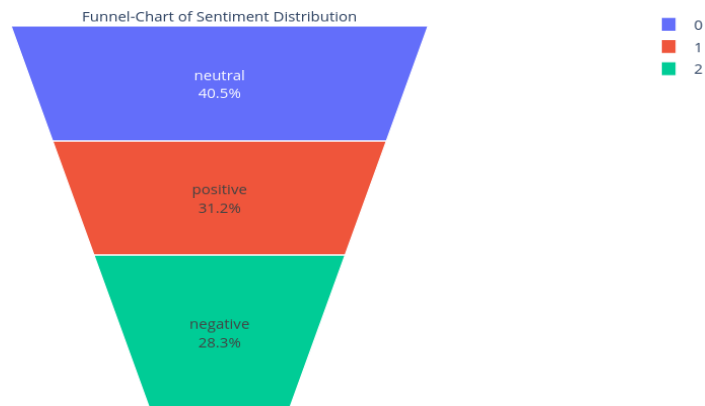
```
1 plt.figure(figsize=(12,6))
2 sns.countplot(x='sentiment',data=train)
3
```

```
Out[11]: <matplotlib.axes._subplots.AxesSubplot at 0x7f99b80710b8>
```



برای مصورسازی بهتر، یک Funnel-chart از داده ترسیم می‌کنیم:

```
1 fig = go.Figure(go.Funnelarea(
2     text=temp.sentiment,
3     values=temp.values,
4     title={"position": "top center", "text": "Funnel-Chart of
5         Sentiment Distribution"}
6 ))
7 fig.show()
```



مواردی که تا کنون از داده می‌دانیم:

- می‌دانیم `selected_text` زیرمجموعه‌یی از `text` است.
- می‌دانیم `selected_text` تنها یک بخش از `text` را شامل می‌شود. یعنی، از بخشی از یک عبارت و بخشی از یک عبارت دیگر را شامل نمی‌شود. برای مثال؛ اگر `text` برابر با:

«Spent the entire morning in a meeting w/ a vendor, and my boss was not happy w/ them. Lots of fun. I had other plans for my morning»

باشد؛ selected_text می‌تواند «my boss was not happy w/ them. Lots of fun» یا my «Morning vendor, I had other plans for my morning» اما نمی‌تواند برابر با boss was not happy w/ them. Lots of fun and my boss باشد.

○ با وجود گفتگوی کگلی که از طریق این لینک قابل دسترسی است، می‌دانیم که توپیت‌های بی‌طرف، دارای ۹۷ درصد شباهت ژاکارد بین text و selected_text می‌باشد.

○ همچنین با استفاده از گفتگویی که از طریق این لینک در اختیارمان قرار دارد، می‌دانیم selected_text در بعضی سطرها از میان یک کلمه شروع می‌شود، در نتیجه selected_text همیشه معنادار نیست. از آنجایی که ما نمی‌دانیم خروجی مجموعه تست هم این ناسازگاری را دارد یا نه، مطمئن نیستیم حذف علائم نگارشی رویکرد خوبی باشد.

تولید متاچرها

از آنجایی که می‌خواهیم selected_text-که زیرمجموعه‌یی از text است-را پیش‌بینی کنیم، ویژگی‌هایی که تولیدشان سودمند خواهد بود عبارت است از:

○ اختلاف بین تعداد کلمات در text و selected_text

○ ضریب شباهت ژاکارد برای text و selected_text

تابع برای ضریب ژاکارد:

```
1 def jaccard(str1, str2):
2     a = set(str1.lower().split())
3     b = set(str2.lower().split())
4     c = a.intersection(b)
5     return float(len(c)) / (len(a) + len(b) - len(c))
```

به‌دست آوردن ضریب ژاکارد برای تمام سطرهای داده‌ی آموزشی:

```
1 results_jaccard=[]
2
3 for ind,row in train.iterrows():
4     sentence1 = row.text
5     sentence2 = row.selected_text
6
7     jaccard_score = jaccard(sentence1,sentence2)
8     results_jaccard.append([sentence1,sentence2, jaccard_score])
```

تبدیل لیست ضریب ژاکارد-که در سلول قبلی محاسبه شد-به دیتافریم و ادغام آن با داده‌ی آموزشی:

```
1 jaccard = pd.DataFrame(results_jaccard,columns=["text","selected_text",
2     "jaccard_score"])
3 train = train.merge(jaccard,how='outer')
```

اضافه کردن ویژگی‌های موردنظر - Num_words_ST, Num_word_text, و difference_in_words - به داده‌ی آموزشی

```
1 train['Num_words_ST'] = train['selected_text'].apply(lambda x:len(str(x)
2 ).split())) #Number Of words in Selected Text
3 train['Num_word_text'] = train['text'].apply(lambda x:len(str(x).split
4 ())) #Number Of words in main text
5 train['difference_in_words'] = train['Num_word_text'] - train['
6 Num_words_ST'] #Difference in Number of words text and Selected
7 Text
```

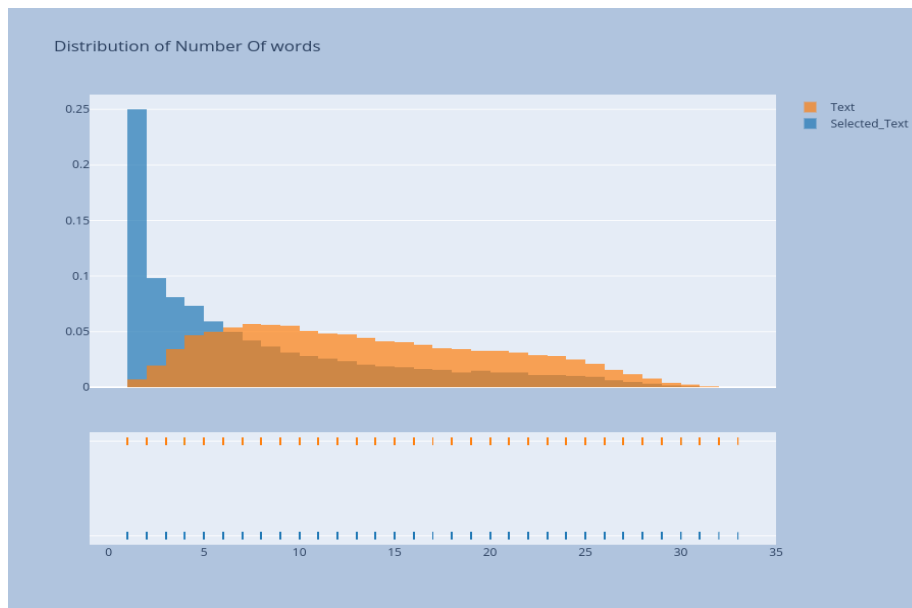
نمونه‌یی از داده‌ی آموزشی بعد از اضافه کردن متافچرها:

Out[17]:

| | textID | text | selected_text | sentiment | jaccard_score | Num_words_ST | Num_word_text | difference_in_words |
|---|------------|--------------------------------------------------|-------------------------------------|-----------|---------------|--------------|---------------|---------------------|
| 0 | cb774db0d1 | I'd have responded, if I were going | I'd have responded, if I were going | neutral | 1.000000 | 7 | 7 | 0 |
| 1 | 549e992a42 | Sooo SAD I will miss you here in San Diego!!! | Sooo SAD | negative | 0.200000 | 2 | 10 | 8 |
| 2 | 088c60f138 | my boss is bullying me... | bullying me | negative | 0.166667 | 2 | 5 | 3 |
| 3 | 9642c003ef | what interview! leave me alone | leave me alone | negative | 0.600000 | 3 | 5 | 2 |
| 4 | 358bd9e861 | Sons of ****, why couldn't they put them on t... | Sons of ****, | negative | 0.214286 | 3 | 14 | 11 |

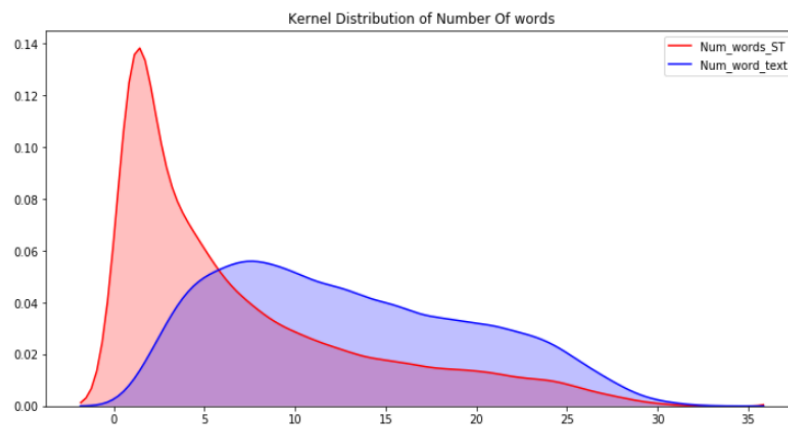
حال بیا باید نگاهی به توزیع متافچرها بی‌اندازیم.

```
1 hist_data = [train['Num_words_ST'],train['Num_word_text']]
2
3 group_labels = ['Selected_Text', 'Text']
4
5 # Create distplot with custom bin_size
6 fig = ff.create_distplot(hist_data, group_labels,show_curve=False)
7 fig.update_layout(title_text='Distribution of Number Of words')
8 fig.update_layout(
9     autosize=False,
10     width=900,
11     height=700,
12     paper_bgcolor="LightSteelBlue",
13 )
14 fig.show()
15
```



نمودار تعداد کلمات از این نظر حائز اهمیت است که توییت‌هایی که تعداد کلمات‌شان بیشتر از ۲۵ تا می‌باشد بسیار کم‌تر هستند، از این رو منحنی توزیع چوله به راست است.

```
1 p1=sns.kdeplot(train['Num_words_ST'], shade=True, color="r").set_title('Kernel Distribution of Number Of words')
2 p1=sns.kdeplot(train['Num_word_text'], shade=True, color="b")
```

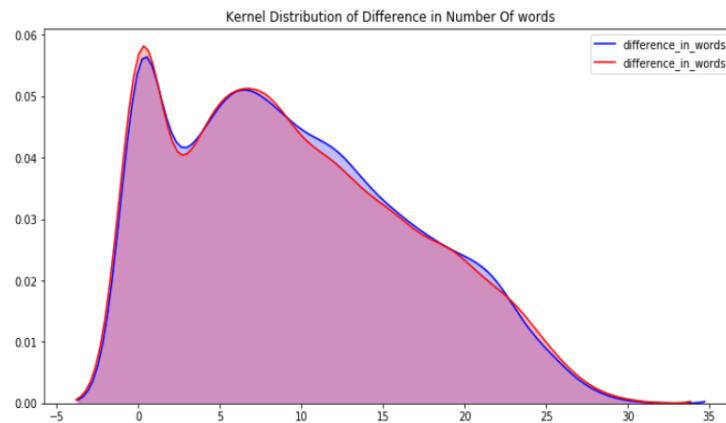


حالا مفید است نگاهی هم به توزیع اختلاف تعداد کلمات و ضریب ژاکارد در تمایل‌ها بی‌اندازیم.

```

1 plt.figure(figsize=(12,6))
2 p1=sns.kdeplot(train[train['sentiment']=='positive']['
    difference_in_words'], shade=True, color="b").set_title('Kernel
    Distribution of Difference in Number Of words')
3 p2=sns.kdeplot(train[train['sentiment']=='negative']['
    difference_in_words'], shade=True, color="r")
4

```



```

1 plt.figure(figsize=(12,6))
2 sns.distplot(train[train['sentiment']=='neutral']['difference_in_words'
    ],kde=False)
3

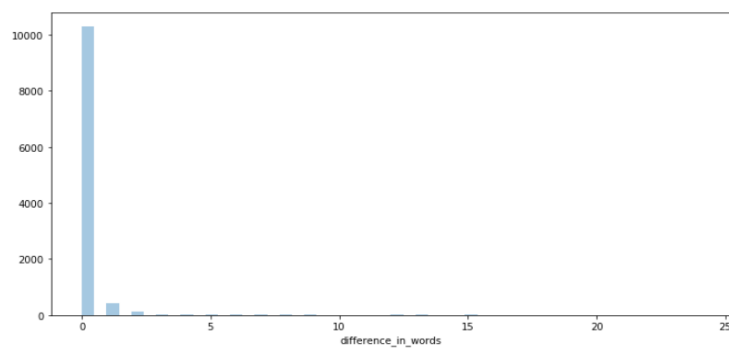
```

Out[21]:

```

<matplotlib.axes._subplots.AxesSubplot at 0x7f9968f22ac8>

```



قادر به ترسیم KDE Plot برای توییت‌های بی‌تمایل (Neutral) نیستیم، چراکه برای اکثر توییت‌های خنثی اختلاف میان تعداد کلمات در text و selected_text برابر صفر است.

```

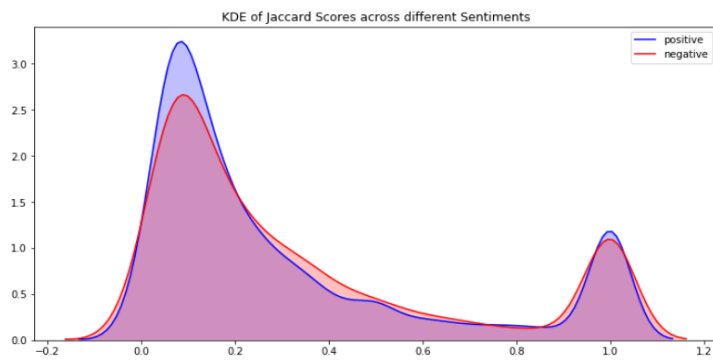
1 plt.figure(figsize=(12,6))
2 p1=sns.kdeplot(train[train['sentiment']=='positive']['jaccard_score'],
3               shade=True, color="b").set_title('KDE of Jaccard Scores across
4               different Sentiments')
5 p2=sns.kdeplot(train[train['sentiment']=='negative']['jaccard_score'],
6               shade=True, color="r")
7 plt.legend(labels=['positive','negative'])

```

```

Out[22]:
<matplotlib.legend.Legend at 0x7f99682eac50>

```



به همان دلیل ذکر شده در مورد ترسیم KDE Plot برای اختلاف کلمات در توپیت‌های خنثی، قادر به ترسیم KDE Plot برای ضرب زاکارد در توپیت‌های خنثی نیز نیستیم.

```

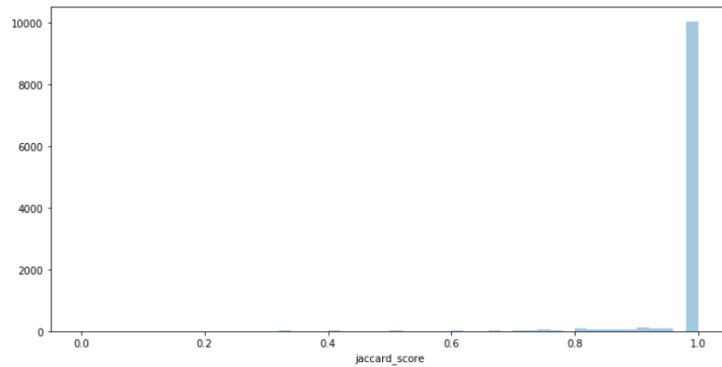
1 plt.figure(figsize=(12,6))
2 sns.distplot(train[train['sentiment']=='neutral']['jaccard_score'],kde=
3               False)

```

```

Out[23]:
<matplotlib.axes._subplots.AxesSubplot at 0x7f99682fce80>

```



ما می‌توانیم چند گرایش جالب را در اینجا ببینیم:

- توییت‌های مثبت و منفی کورتوز بالایی دارند، از این رو، توییت‌ها در دو ناحیه‌ی کم‌تراکم و پر تراکم متمرکز شده‌اند.
- توییت‌های خنثی مقدار کورتوز پایینی دارند و یک برآمدگی تراکم در نزدیکی مقادیر ۱ وجود دارد.

نتیجه‌گیری

- با مشاهده‌ی پلات ضریب ژاکارد، می‌توان دید برای توییت‌های مثبت و منفی، در نزدیکی ۱ برآمدگی وجود دارد. این یعنی، یک خوشه از توییت‌ها وجود دارد که در آن‌ها text و selected_text شباهت بالایی دارند. اگر بتوانیم آن خوشه‌ها را پیدا کنیم، با صرف نظر از کلاس‌های sentiment، می‌توانیم selected_text را پیش‌بینی کنیم.

بیا یاد بررسی کنیم که آیا می‌توان این خوشه‌ها را پیدا کرد:

```
1 k = train[train['Num_word_text']<=2]
2 k.groupby('sentiment').mean()['jaccard_score']
```

```
Out[25]:
sentiment
negative    0.788580
neutral     0.977805
positive    0.765700
Name: jaccard_score, dtype: float64
```

می‌توان دید شباهتی میان text و selected_text وجود دارد. بیا یاد با دقت بیشتر بررسی کنیم:

```
1 k[k['sentiment']=='positive']
2
```

واضح است که بیشتر مواقع، text و selected_text برابر هستند.

پاک‌سازی داده

حالا استخراج اطلاعات از کلمات توییت‌ها را آغاز خواهیم کرد. بیا یاد اول داده‌ی آموزشی را پاک‌سازی کنیم:

```
1 def clean_text(text):
2     '''Make text lowercase, remove text in square brackets,remove links,
3     remove punctuation
4     and remove words containing numbers.'''
5     text = str(text).lower()
6     text = re.sub('\[.*?\]', '', text)
7     text = re.sub('https?://\S+|www.\S+', '', text)
8     text = re.sub('<.*?>+', '', text)
9     text = re.sub('[%s]' % re.escape(string.punctuation), '', text)
10    text = re.sub('\n', '', text)
```

Out[26]:

| | textID | text | selected_text | sentiment | jaccard_score | Num |
|-------|------------|------------------------------------|------------------------------------|-----------|---------------|-----|
| 68 | fa2654e730 | Chilllin | Chilllin | positive | 1.0 | 1 |
| 80 | bbbc46889b | THANK YYYYYYYYYOOOOOOOOOOUUUUUI | THANK YYYYYYYYYOOOOOOOOOOUUUUUI | positive | 1.0 | 2 |
| 170 | f3d95b57b1 | good morning | good morning | positive | 1.0 | 2 |
| 278 | 89d5b3f0b5 | Thanks | Thanks | positive | 1.0 | 1 |
| 429 | a78ef3e0d0 | Goodmorning | Goodmorning | positive | 1.0 | 1 |
| ... | ... | ... | ... | ... | ... | ... |
| 26689 | e80c242d6a | Goodnight; | Goodnight; | positive | 1.0 | 1 |
| 26725 | aad244f37d | *hug* | *hug* | positive | 1.0 | 1 |
| 26842 | a46571fe12 | congrats! | congrats! | positive | 1.0 | 1 |
| 26959 | 49a942e9b1 | Happy birthday. | Happy birthday. | positive | 1.0 | 2 |
| 27292 | 47c474aaf1 | Good choice | Good | positive | 0.5 | 1 |

```
10 text = re.sub('\w*\d\w*', '', text)
11 return text
12
```

```
1 train['text'] = train['text'].apply(lambda x:clean_text(x))
2 train['selected_text'] = train['selected_text'].apply(lambda x:
3 clean_text(x))
```

```
1 train.head()
2
```

Out[29]:

| | textID | text | selected_text | sentiment | jaccard_score | Num_words_ST | Num_word_text | difference_in_words |
|---|-------------|----------------------------------------------------------------|-----------------------------------------|-----------|---------------|--------------|---------------|---------------------|
| 0 | cb7774db0d1 | id have responded if i were going | id have responded if i were going | neutral | 1.000000 | 7 | 7 | 0 |
| 1 | 549e992a42 | sooo sad i will miss you here in san diego | sooo sad | negative | 0.200000 | 2 | 10 | 8 |
| 2 | 088c60f138 | my boss is bullying me | bullying me | negative | 0.166667 | 2 | 5 | 3 |
| 3 | 9642c003ef | what interview leave me alone | leave me alone | negative | 0.600000 | 3 | 5 | 2 |
| 4 | 358bd9e861 | sons of why couldnt they put them on the rel... | sons of | negative | 0.214286 | 3 | 14 | 11 |

رایج ترین کلمات در **selected_text**

```
1 train['temp_list'] = train['selected_text'].apply(lambda x:str(x).split
    ())
2 top = Counter([item for sublist in train['temp_list'] for item in
    sublist])
3 temp = pd.DataFrame(top.most_common(20))
4 temp.columns = ['Common_words', 'count']
5 temp.style.background_gradient(cmap='Blues')
6
```


Out[30]:

| | Common_words | count |
|----|--------------|-------|
| 0 | i | 7200 |
| 1 | to | 5305 |
| 2 | the | 4590 |
| 3 | a | 3538 |
| 4 | my | 2783 |
| 5 | you | 2624 |
| 6 | and | 2321 |
| 7 | it | 2158 |
| 8 | is | 2115 |
| 9 | in | 1986 |
| 10 | for | 1854 |
| 11 | im | 1676 |
| 12 | of | 1638 |
| 13 | me | 1540 |
| 14 | on | 1488 |
| 15 | so | 1410 |
| 16 | have | 1345 |
| 17 | that | 1297 |
| 18 | but | 1267 |
| 19 | good | 1251 |

مشاهده می‌کنید که stopword ها را حذف نکرده‌ایم (برای نمونه: کلمه ی to).

```
1 def remove_stopword(x):  
2     return [y for y in x if y not in stopwords.words('english')]  
3 train['temp_list'] = train['temp_list'].apply(lambda x:remove_stopword(  
    x))
```

Out[33]:

| | Common_words | count |
|----|--------------|-------|
| 1 | good | 1251 |
| 2 | day | 1058 |
| 3 | love | 909 |
| 4 | happy | 852 |
| 5 | like | 774 |
| 6 | get | 772 |
| 7 | dont | 765 |
| 8 | go | 700 |
| 9 | cant | 613 |
| 10 | work | 612 |
| 11 | going | 592 |
| 12 | today | 564 |
| 13 | got | 558 |
| 14 | one | 538 |
| 15 | time | 534 |
| 16 | thanks | 532 |
| 17 | lol | 528 |
| 18 | really | 520 |
| 19 | u | 519 |

رایج‌ترین کلمات در **text**

```

1 train['temp_list1'] = train['text'].apply(lambda x:str(x).split()) #
    List of words in every row for text
2 train['temp_list1'] = train['temp_list1'].apply(lambda x:
    remove_stopword(x)) #Removing Stopwords

1 top = Counter([item for sublist in train['temp_list1'] for item in
    sublist])
2 temp = pd.DataFrame(top.most_common(25))
3 temp = temp.iloc[1:,:]
4 temp.columns = ['Common_words','count']
5 temp.style.background_gradient(cmap='Blues')

```

Out[36]:

| | Common_words | count |
|----|--------------|-------|
| 1 | day | 2044 |
| 2 | good | 1549 |
| 3 | get | 1426 |
| 4 | like | 1346 |
| 5 | go | 1267 |
| 6 | dont | 1200 |
| 7 | love | 1122 |
| 8 | work | 1112 |
| 9 | going | 1096 |
| 10 | today | 1096 |
| 11 | got | 1072 |
| 12 | cant | 1020 |
| 13 | happy | 976 |
| 14 | one | 971 |
| 15 | lol | 948 |
| 16 | time | 942 |
| 17 | know | 930 |
| 18 | u | 923 |
| 19 | really | 908 |
| 20 | back | 891 |
| 21 | see | 797 |
| 22 | well | 744 |
| 23 | new | 740 |
| 24 | night | 737 |

رایج ترین کلمات در تمایل ها

```

1 Positive_sent = train[train['sentiment']=='positive']
2 Negative_sent = train[train['sentiment']=='negative']
3 Neutral_sent = train[train['sentiment']=='neutral']

```

رایج ترین کلمات در توییت های مثبت

```

1 #Most common positive words
2 top = Counter([item for sublist in Positive_sent['temp_list'] for item
    in sublist])
3 temp_positive = pd.DataFrame(top.most_common(20))
4 temp_positive.columns = ['Common_words','count']
5 temp_positive.style.background_gradient(cmap='Greens')

```

Out[39]:

| | Common_words | count |
|----|--------------|-------|
| 0 | good | 826 |
| 1 | happy | 730 |
| 2 | love | 697 |
| 3 | day | 456 |
| 4 | thanks | 439 |
| 5 | great | 364 |
| 6 | fun | 287 |
| 7 | nice | 267 |
| 8 | mothers | 259 |
| 9 | hope | 245 |
| 10 | awesome | 232 |
| 11 | im | 185 |
| 12 | thank | 180 |
| 13 | like | 167 |
| 14 | best | 154 |
| 15 | wish | 152 |
| 16 | amazing | 135 |
| 17 | really | 128 |
| 18 | better | 125 |
| 19 | cool | 119 |

رایج ترین کلمات در توییت های منفی

```
1 #Most common negative words
2 top = Counter([item for sublist in Negative_sent['temp_list'] for item
3               in sublist])
4 temp_negative = pd.DataFrame(top.most_common(20))
5 temp_negative = temp_negative.iloc[1:,:]
6 temp_negative.columns = ['Common_words', 'count']
7 temp_negative.style.background_gradient(cmap='Reds')
```

Out[41]:

| | Common_words | count |
|----|--------------|-------|
| 1 | miss | 358 |
| 2 | sad | 343 |
| 3 | sorry | 300 |
| 4 | bad | 246 |
| 5 | hate | 230 |
| 6 | dont | 221 |
| 7 | cant | 201 |
| 8 | sick | 166 |
| 9 | like | 162 |
| 10 | sucks | 159 |
| 11 | feel | 158 |
| 12 | tired | 144 |
| 13 | really | 137 |
| 14 | good | 127 |
| 15 | bored | 115 |
| 16 | day | 110 |
| 17 | hurts | 108 |
| 18 | work | 99 |
| 19 | get | 97 |

رایج ترین کلمات در توییت های خنثی

```
1 #Most common Neutral words
2 top = Counter([item for sublist in Neutral_sent['temp_list'] for item
3               in sublist])
4 temp_neutral = pd.DataFrame(top.most_common(20))
5 temp_neutral = temp_neutral.loc[1:,:]
6 temp_neutral.columns = ['Common_words', 'count']
7 temp_neutral.style.background_gradient(cmap='Reds')
```

Out[41]:

| | Common_words | count |
|----|--------------|-------|
| 1 | miss | 358 |
| 2 | sad | 343 |
| 3 | sorry | 300 |
| 4 | bad | 246 |
| 5 | hate | 230 |
| 6 | dont | 221 |
| 7 | cant | 201 |
| 8 | sick | 166 |
| 9 | like | 162 |
| 10 | sucks | 159 |
| 11 | feel | 158 |
| 12 | tired | 144 |
| 13 | really | 137 |
| 14 | good | 127 |
| 15 | bored | 115 |
| 16 | day | 110 |
| 17 | hurts | 108 |
| 18 | work | 99 |
| 19 | get | 97 |

○ مشاهده می‌شود کلماتی همچون get, go, dont, u, got, cant, lol و like در هر سه کلاس تمایل مشترک هستند، نکته‌ی جالب اینجاست که کلماتی همچون: cant و dont طبیعتاً منفی هستند و کلمه‌یی همچون: lol، طبیعتاً مثبت است. آیا داده‌ی ما اشتباه برچسب خورده؟-بعد از آنالیز N-gram دید بهتری به این موضوع خواهیم داشت.

○ بررسی کلمات یکتا در تمایل‌های مختلف جالب است.

کلمات یکتا در هر تمایل

ما کلمات متمایز را برای هر تمایل-با استفاده از تابع words_unique-به ترتیب زیر بررسی خواهیم کرد:

○ مثبت

○ منفی

○ خنثی

```
1 raw_text = [word for word_list in train['temp_list1'] for word in
2   word_list]
```

توییت‌های مثبت

```
1 Unique_Positive= words_unique('positive', 20, raw_text)
2 print("The top 20 unique words in Positive Tweets are:")
3 Unique_Positive.style.background_gradient(cmap='Greens')
4
```

The top 20 unique words in Positive Tweets are:

Out[48]:

| | words | count |
|----|-----------------|-------|
| 0 | congratulations | 29 |
| 1 | thnx | 10 |
| 2 | appreciated | 8 |
| 3 | shared | 7 |
| 4 | presents | 7 |
| 5 | greetings | 7 |
| 6 | blessings | 6 |
| 7 | mothersday | 6 |
| 8 | mcr | 6 |
| 9 | coolest | 6 |
| 10 | honored | 6 |
| 11 | goood | 6 |
| 12 | wango | 5 |
| 13 | actress | 5 |
| 14 | mint | 5 |
| 15 | dayyyy | 5 |
| 16 | clara | 5 |
| 17 | twin | 5 |
| 18 | kudos | 5 |
| 19 | hurray | 5 |

```

1 fig = px.treemap(Unique_Positive, path=['words'], values='count',title=
  'Tree Of Unique Positive Words')
2 fig.show()
3

```

Tree Of Unique Positive Words



```

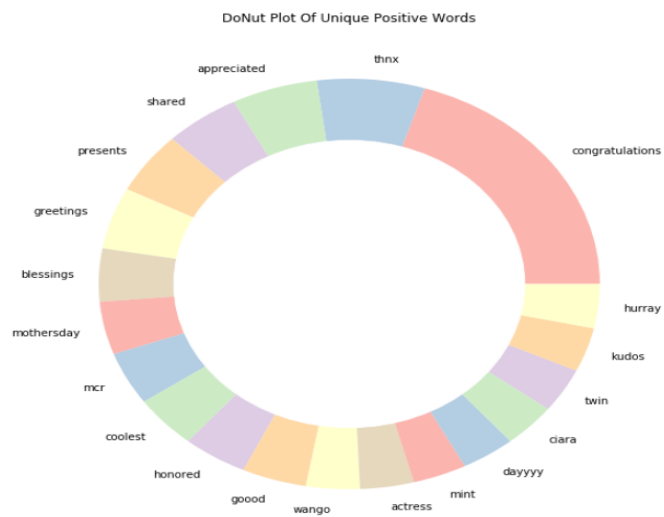
1 from palettable.colorbrewer.qualitative import Pastel1_7
2 plt.figure(figsize=(16,10))
3 my_circle=plt.Circle((0,0), 0.7, color='white')
4 plt.pie(Unique_Positive['count'], labels=Unique_Positive.words, colors=
  Pastel1_7.hex_colors)
5 p=plt.gcf()
6 p.gca().add_artist(my_circle)

```

```

7 plt.title('DoNut Plot Of Unique Positive Words')
8 plt.show()
9

```



توییت‌های منفی

```

1 Unique_Negative= words_unique('negative', 10, raw_text)
2 print("The top 10 unique words in Negative Tweets are:")
3 Unique_Negative.style.background_gradient(cmap='Reds')
4

```

The top 20 unique words in Positive Tweets are:

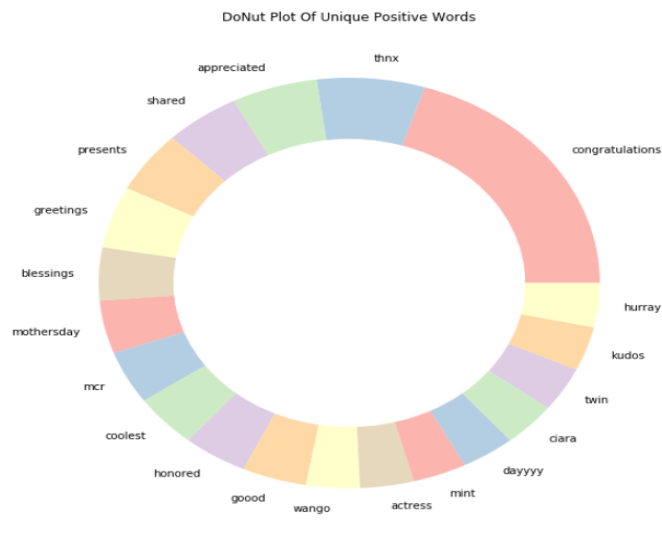
out[48]:

| | words | count |
|----|-----------------|-------|
| 0 | congratulations | 29 |
| 1 | thnx | 10 |
| 2 | appreciated | 8 |
| 3 | shared | 7 |
| 4 | presents | 7 |
| 5 | greetings | 7 |
| 6 | blessings | 6 |
| 7 | mothersday | 6 |
| 8 | mcr | 6 |
| 9 | coolest | 6 |
| 10 | honored | 6 |
| 11 | goood | 6 |
| 12 | wango | 5 |
| 13 | actress | 5 |
| 14 | mint | 5 |
| 15 | dayyyy | 5 |
| 16 | clara | 5 |
| 17 | twin | 5 |
| 18 | kudos | 5 |
| 19 | hurray | 5 |

```

1 from palettable.colorbrewer.qualitative import Pastel1_7
2 plt.figure(figsize=(16,10))
3 my_circle=plt.Circle((0,0), 0.7, color='white')
4 plt.rcParams['text.color'] = 'black'
5 plt.pie(Unique_Negative['count'], labels=Unique_Negative.words, colors=
6         Pastel1_7.hex_colors)
7 p=plt.gcf()
8 p.gca().add_artist(my_circle)
9 plt.title('DoNut Plot Of Unique Negative Words')
10 plt.show()

```



توییت‌های خنثی

```

1 Unique_Neutral= words_unique('neutral', 10, raw_text)
2 print("The top 10 unique words in Neutral Tweets are:")
3 Unique_Neutral.style.background_gradient(cmap='Oranges')
4

```

The top 20 unique words in Positive Tweets are:

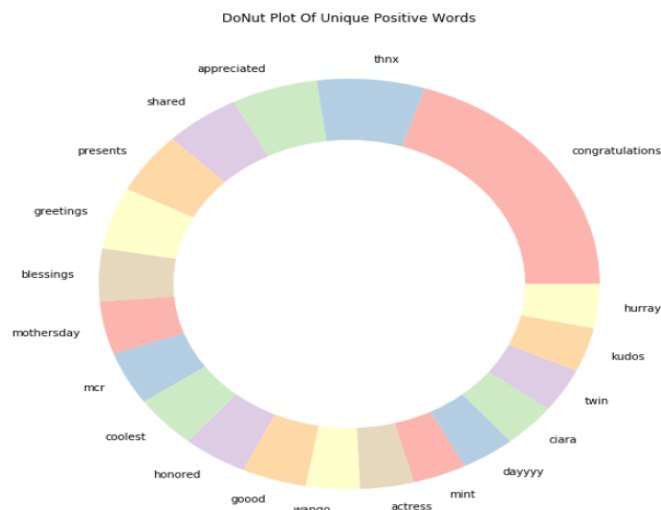
out[48]:

| | words | count |
|----|-----------------|-------|
| 0 | congratulations | 29 |
| 1 | thnx | 10 |
| 2 | appreciated | 8 |
| 3 | shared | 7 |
| 4 | presents | 7 |
| 5 | greetings | 7 |
| 6 | blessings | 6 |
| 7 | mothersday | 6 |
| 8 | mcr | 6 |
| 9 | coolest | 6 |
| 10 | honored | 6 |
| 11 | goood | 6 |
| 12 | wango | 5 |
| 13 | actress | 5 |
| 14 | mint | 5 |
| 15 | dayyyy | 5 |
| 16 | clara | 5 |
| 17 | twin | 5 |
| 18 | kudos | 5 |
| 19 | hurray | 5 |

```

1 from palettable.colorbrewer.qualitative import Pastel1_7
2 plt.figure(figsize=(16,10))
3 my_circle=plt.Circle((0,0), 0.7, color='white')
4 plt.pie(Unique_Neutral['count'], labels=Unique_Neutral.words, colors=
5         Pastel1_7.hex_colors)
6 p=plt.gcf()
7 p.gca().add_artist(my_circle)
8 plt.title('DoNut Plot Of Unique Neutral Words')
9 plt.show()

```



حال که کلمات متمایز را بررسی کردیم، دید به مراتب واضح‌تری از داده داریم. این کلمات یکتا از تعیین کنندگان قوی sentiment برای توییت‌ها می‌باشند.

مدل‌سازی

این مسئله می‌تواند به شکل‌های زیر مدل شود:

○ Named Entity Recognition

○ Question-Answering

○ من یک رویکرد دیگر-که Nick در کرنل‌اش به اشتراک گذاشته-یافتیم که از مفهوم Gini Impurity استفاده کرده تا به کلمات درون توییت‌ها وزن بدهد، سپس پیش‌بینی بر اساس آن وزن‌ها اتفاق می‌افتد. می‌توانید از طریق این لینک به کرنل رجوع کنید.

○ یک ایده‌ی متفاوت‌تر از طریق این لینک قابل دسترسی است.

○ یک رویکرد مفید دیگر: لینک

منابع

○ برای مدلینگ مسئله به‌عنوان مسئله‌ی NER: لینک

○ برای مدلینگ مسئله به‌عنوان مسئله‌ی Question-Answering: لینک

مدل کردن مسئله به‌عنوان مسئله‌ی NER

NER یک استاندارد برای مسئله‌های پردازش زبان طبیعی است، که شامل استخراج موجودیت‌های اسم‌دار-مردم، مکان‌ها، سازمان‌ها و غیره-از داده‌ی متنی و طبقه‌بندی آن‌ها در مجموعه‌ای از دسته‌ها می‌باشد. برای درک بهتر NER به این لینک مراجعه کنید. مواردی در مورد مدل‌مان:

- ما برای توییت‌های بی‌تمایل، از text به جای selected_text استفاده خواهیم کرد-به‌دلیل ضریب بالای ژاکارد در آن‌ها.
- همچنین در توییت‌هایی که کم‌تر از سه کلمه دارند هم text را به جای selected_text به‌کار خواهیم برد.
- دو مدل را برای توییت‌های مثبت و توییت‌های منفی آموزش داده می‌شوند.
- داده پیش‌پردازش نخواهد شد، چرا که selected_text تنها شامل متن خام می‌شود.

خواندن داده‌های آموزشی و آزمایشی

```
1 df_train = pd.read_csv('/kaggle/input/tweet-sentiment-extraction/train.csv')
2 df_test = pd.read_csv('/kaggle/input/tweet-sentiment-extraction/test.csv')
3 df_submission = pd.read_csv('/kaggle/input/tweet-sentiment-extraction/sample_submission.csv')
```

رایج‌ترین کلمات در توییت‌های منفی

```
1 df_train['Num_words_text'] = df_train['text'].apply(lambda x:len(str(x).split())) #Number Of words in main Text in train set
```

انتخاب توییت‌های که بیشتر از دو کلمه را شامل می‌شوند

```
1 df_train = df_train[df_train['Num_words_text']>=3]
```

تابعی برای ذخیره کردن مدل در مسیر داده شده

```
1 def save_model(output_dir, nlp, new_model_name):
2     ''' This Function Saves model to
3     given output directory'''
4
5     output_dir = f'../working/{output_dir}'
6     if output_dir is not None:
7         if not os.path.exists(output_dir):
8             os.makedirs(output_dir)
9             nlp.meta["name"] = new_model_name
10            nlp.to_disk(output_dir)
11            print("Saved model to", output_dir)
12
```

تابعی برای دریافت کردن مسیر مدل‌ها

```
1 def get_model_out_path(sentiment):
2     '''
3     Returns Model output path
```

```

4     '''
5     model_out_path = None
6     if sentiment == 'positive':
7         model_out_path = 'models/model_pos'
8     elif sentiment == 'negative':
9         model_out_path = 'models/model_neg'
10    return model_out_path
11

```

تابعی برای دریافت کردن داده با فرمتی که برای **Spacy NER** مناسب است

```

1 def get_training_data(sentiment):
2     '''
3     Returns Trainong data in the format needed to train spacy NER
4     '''
5     train_data = []
6     for index, row in df_train.iterrows():
7         if row.sentiment == sentiment:
8             selected_text = row.selected_text
9             text = row.text
10            start = text.find(selected_text)
11            end = start + len(selected_text)
12            train_data.append((text, {"entities": [[start, end, 'selected_text'
13                ]]}))
14    return train_data

```

آموزش مدل‌ها

```

1 sentiment = 'positive'
2
3 train_data = get_training_data(sentiment)
4 model_path = get_model_out_path(sentiment)
5 # For DEmo Purposes I have taken 3 iterations you can train the model
6   as you want
7 train(train_data, model_path, n_iter=3, model=None)

```

```

1 sentiment = 'negative'
2
3 train_data = get_training_data(sentiment)
4 model_path = get_model_out_path(sentiment)
5
6 train(train_data, model_path, n_iter=3, model=None)
7

```

پیشبینی با مدل‌های آموزش داده شده

```

1 def predict_entities(text, model):
2     doc = model(text)
3     ent_array = []

```

```

4     for ent in doc.ents:
5         start = text.find(ent.text)
6         end = start + len(ent.text)
7         new_int = [start, end, ent.label_]
8         if new_int not in ent_array:
9             ent_array.append([start, end, ent.label_])
10        selected_text = text[ent_array[0][0]: ent_array[0][1]] if len(
11            ent_array) > 0 else text
12        return selected_text

```

```

1  MODELS_BASE_PATH = '../input/tse-spacy-model/models/'
2
3  if MODELS_BASE_PATH is not None:
4      print("Loading Models from ", MODELS_BASE_PATH)
5      model_pos = spacy.load(MODELS_BASE_PATH + 'model_pos')
6      model_neg = spacy.load(MODELS_BASE_PATH + 'model_neg')
7
8      for index, row in df_test.iterrows():
9          text = row.text
10         output_str = ""
11         if row.sentiment == 'neutral' or len(text.split()) <= 2:
12             selected_texts.append(text)
13         elif row.sentiment == 'positive':
14             selected_texts.append(predict_entities(text, model_pos))
15         else:
16             selected_texts.append(predict_entities(text, model_neg))
17     df_test['selected_text'] = selected_texts

```

```

Loading Models from ../input/tse-spacy-model/models/

```

```

1 df_submission['selected_text'] = df_test['selected_text']
2 df_submission.to_csv("submission.csv", index=False)
3 display(df_submission.head(10))
4

```

| | textID | selected_text |
|---|------------|--------------------------------------------------|
| 0 | f87dea47db | Last session of the day http://twitpic.com/67ezh |
| 1 | 96d74cb729 | exciting |
| 2 | eee518ae67 | Recession |
| 3 | 01082688c6 | happy bday! |
| 4 | 33987a8ee5 | I like it!! |
| 5 | 726e501993 | visitors! |
| 6 | 261932614e | HATES |
| 7 | afa11da83f | blocked |
| 8 | e64208b4ef | and within a short time of the last clue all ... |
| 9 | 37bcad24ca | What did you get? My day is alright.. haven`... |

سخن پایانی

کگل همیشه فرصت زیادی برای رقابت‌ها در نظر می‌گیرد، اینگونه شرکت‌کنندگان می‌توانند یاد بگیرند و رشد کنند. همانطور که عهد کرده بودیم یک مدل همراه با توضیحات‌اش را ارائه کردیم؛ برای مطالعه بیشتر، می‌توانید مستندات `spacy` را در این کرنل مطالعه کنید.

References

- [1] A simple solution using only word counts
- [2] Twitter sentiment Extaction-Analysis, EDA and Model
- [3] NER - training using spacy (Ensemble)
- [4] Question-Answering Starter pack