

اطلاعات کرال شده را در نهایت در یک فایل json ذخیره سازی میکنیم.

```
{
  {
    "news_headline": "Godzilla x Kong: The New Empire review \u2013 breezy, forgettable monster sequel",
    "news_firstline": "There\u2019s a likable, light-hearted zip to the monster mash follow-up but energy dissipates when we\u2019re stuck with the humans",
    "rating": 3,
    "news_article": "It was a strange old time when the creature feature mash-up Godzilla vs Kong was released, the first major blockbuster in cinemas since Covid shuttered them
  },
  {
    "news_headline": "Mary Poppins review \u2013 Disney\u2019s entertainment sugar rush possesses thermonuclear brilliance",
    "news_firstline": "Manic, magic, madcap \u2013 Julie Andrews is superb in the role of the flying nanny, in a film filled with amazing songs",
    "rating": 5,
    "news_article": "Brilliant, entrancing, exhausting, and with thermonuclear showtunes from Richard and Robert Sherman, Disney\u2019s hybrid live-action/animation classic from 1964
  },
  {
    "news_headline": "Kung Fu Panda 4 review \u2013 Jack Black and Awkwafina in hurricane of slapstick more miss than hit",
    "news_firstline": "The lead pair make a brilliant double act, but the franchise has run out of its signature sweetness and charm",
    "rating": 2,
    "news_article": "The cuddly kung fu master is back. Jack Black returns as dumpling-loving panda Po, the unlikely of lean, mean fightin\u2019 machines. It\u2019s been eight years since the last film
  },
  {
    "news_headline": "Ryuichi Sakamoto: Opus review \u2013 a stark, emotional finale from master musician",
    "news_firstline": "In his last weeks of life, the Oscar-winning composer is filmed at the piano by his son. It is an almost wordless paean to a remarkable career",
    "rating": 4,
    "news_article": "Short of presenting nothing more than music and a blank screen, this documentary about the late Japanese composer-performer Ryuichi Sakamoto\u2019s last appearance
  },
  {
    "news_headline": "The Persian Version review \u2013 feelgood Iranian-American comedy with edge",
    "news_firstline": "A perceptive and candid look at mother-daughter discord drives the boisterous energy of Maryam Keshavarz\u2019s comic crowdpleaser",
    "rating": 3,
    "news_article": "The only daughter in a family of nine Iranian-American siblings, Leila\u2019s (Layla Mohammadi) relationship with her mother, Shireen (Niousha Noor), was set against the backdrop of the 2015 election
  },
  {
    "news_headline": "Robot Dreams review \u2013 bittersweet buddy movie is one of the best animations in recent years",
    "news_firstline": "A lonely dog buys himself a robot companion and learns to see the world in a joyous new light in Spanish director Pablo Berger\u2019s exquisite, Oscar-nominee
    "rating": 5,
    "news_article": "It\u2019s an almost entirely dialogue-free animation, captured with pleasingly simple, almost naive 2D character design. The warm and disarming storytelling
  },
  {
    "news_headline": "Late Night With the Devil review \u2013 diabolically funny found-footage horror",
    "news_firstline": "A twisted, brilliant, and utterly terrifying film that is a perfect blend of 1970s exploitation and modern horror
  }
}
```

پس از اینکار به سراغ پردازش دادگان میرویم. در بخش اول تسک انتخاب شده به دست آوردن عبارات کلیدی و صفات استفاده شده در این نقد ها میباشد. مراحل مرتبط به کد ها در نوتبوک به شکل مناسب کامنت گذاری شده است.

برای اینکار باید پیش پردازش های گفته شده را انجام دهیم تا انجام اینکار راحت تر شود و اطلاعات با ارزش تر باقی بمانند. در ابتدا تمامی Punctuation ها را حذف میکنیم. سپس تمامی کلمات را Lowercase کرده و در نهایت 's' های موجود در هر کلمه را حذف میکنیم. برای نشان دادن عملکرد از Porter stemmer استفاده کردیم. با استفاده از stemming اتفاقی که می افتد این است که پس از دادن دادگان به Spacy برای شناسایی صفات و عبارات کلیدی این مدل به علت stemming اتفاق افتاده دچار مشکل میشود و نمیتواند این کار را به شکل درست انجام دهد. به عنوان مثال کلمه strange در هنگام stemming به strang تبدیل میشود و آن را verb تشخیص میدهد که این اشتباه با توجه به اینکه تسک ما تشخیص بیشترین صفات به کار رفته در تحلیل فیلم ها میباشد غلط میباشد.

```
1 sentences_a = [sent for sent in doc_preprocessed[0].sents]
2 for x in sentences_a:
3     print([(tok.pos_, tok.lemma_) for tok in x])

✓ [ ('SPACE', ' '), ('VERB', 'strang'), ('ADJ', 'old'), ('NOUN', 'time'), ('NOUN',
'creatur'), ('NOUN', 'featur'), ('NOUN', 'mash'), ('NOUN', 'godzilla'), ('ADP', 'vs'),
('PROPN', 'kong'), ('VERB', 'relea'), ('ADV', 'first'), ('ADJ', 'major'), ('PROPN',
'blockbust'), ('PROPN', 'cinema'), ('PROPN', 'sinc'), ('PROPN', 'covid'), ('PROPN',
```

```
sents = [sent for sent in doc[0].sents]
for x in sents:
    print([(tok.pos_, tok.lemma_) for tok in x])
```

```
[('PRON', 'it'), ('AUX', 'be'), ('DET', 'a'), ('ADJ', 'strange'), ('ADJ', 'old'),
('NOUN', 'time'), ('SCONJ', 'when'), ('DET', 'the'), ('NOUN', 'creature'), ('VERB',
'feature'), ('NOUN', 'mash'), ('PUNCT', '-'), ('NOUN', 'up'), ('PROPN', 'Godzilla'),
('ADP', 'vs'), ('PROPN', 'Kong'), ('AUX', 'be'), ('VERB', 'release'), ('PUNCT', ','),
('DET', 'the'), ('ADJ', 'first'), ('ADJ', 'major'), ('NOUN', 'blockbuster'), ('ADP',
'in'), ('NOUN', 'cinemas'), ('SCONJ', 'since'), ('PROPN', 'Covid'), ('VERB', 'shutter'),
('PRON', 'they'), ('DET', 'all'), ('DET', 'a'), ('NOUN', 'year'), ('ADV', 'prior'),
('PUNCT', '.')]

```

اطلاعات مرتبط با دیتای استفاده شده بعد از پیش پردازش گفته شده در بالا. این پیش پردازش شامل stemming نمیباشد.

```
import itertools
sentence_words = list(itertools.chain(*tokenized_reviews))
print ('%-16s' % 'Number of words', '%-16s' % len(sentence_words))
print ('%-16s' % 'Number of unique words', '%-16s' % len(set(sentence_words)))
avg=np.sum([len(word) for word in sentence_words])/len(sentence_words)
print ('%-16s' % 'Average word length', '%-16s' % avg)
avg_sentence_length_c=np.mean([len(' '.join(sentence)) for sentence in sentence_words])
print ('%-16s' % 'Average sentence length in characters', '%-16s' % avg_sentence_length_c)
```

```
Number of words 327502
Number of unique words 39178
Average word length 6.097312382825143
Average sentence length in characters 11.194624765650286
```

با مقایسه داده ای که این پیش پردازش ها بر روی آن انجام شده در مقابل داده ای که فقط جملات توکنایز شده اند مشخص میشود با توجه به اینکه کلمات حذف شده تاثیری در تسک گفته شده ندارند، با حذف آن ها میتواند اطلاعات را سریعتر و به صورت موثرتری بدست آورد.

همانطور که در تصویر زیر مشخص است توانسته ایم تمامی Noun_chunks ها را بدست آوریم.

```
keywords = []  
for sentence in doc:  
    temp=[]  
    for chunk in sentence.noun_chunks:  
        if chunk.text.lower() not in nlp.Defaults.stop_words:  
            temp.append(chunk.text)  
    keywords.append(temp)
```

```
np.array(keywords)
```

C:\Users\smrh1\AppData\Local\Temp\ipykernel_34188\2185074432.py:1: VisibleDeprecationWarning: Creating an ndarray from nested list of arrays (which results in an ndarray with different lengths or shapes) is deprecated. If you meant to do this, you must specify 'dtype'.
np.array(keywords)

	0
3	[music, a blank screen, this documentar...
4	[The only daughter, a family, nine Iran...
5	[an almost entirely dialogue-free anima...
6	[A malign presence, the corner, poison,...
7	[all the hot takes, the cinema, the sug...
8	[no leprechauns. this abysmal romantic ...

پس از آن بیشترین صفات استفاده شده توسط نویسندگان در این نقد ها را بدست آوردیم.

```

adj_freqdist = FreqDist(itertools.chain(*adjectives))
top40words=adj_freqdist.most_common(40)      # show the top 20 (word, frequency) pairs
print ('%-16s' % 'word', '%-16s' % 'Frequency', '%-16s' % '% of the total')
for topword in top40words:
    percent=(topword[1]/len(tokenized_reviews))
    print ('%-16s' % topword[0], '%-16s' % topword[1], '%-16s' % percent)

```

word	Frequency	% of the total
old	613	0.5125418060200669
new	586	0.4899665551839465
good	561	0.4690635451505017
own	549	0.4590301003344482
young	544	0.45484949832775917
more	538	0.4498327759197324
other	513	0.4289297658862876
first	509	0.42558528428093645
-	491	0.4105351170568562
little	441	0.36872909698996653

در بخش دوم با بررسی هدلاین های هر نقد با توجه به پترن موجود در آن سعی میکنیم که اسم فیلم را استخراج کنیم.

با توجه به اینکه به دنبال اسم فیلم هستیم پیش پردازش خاصی را مد نظر نداریم چون در صورت حذف کردن حتی stopwords اسم فیلم ممکن است متفاوت شود. تنها کاری که میکنیم توکنایز کردن جملات است. سپس با استفاده از regex پترن مدنظر را وارد کرده و سعی میکنیم اسم فیلم را بدست آوریم.

همانطور که در تصویر زیر مشخص است با استفاده از ابزار پردازش متن توانسته ایم با دقت مناسبی اسامی فیلم ها را از هدلاین نقد ها استخراج کنیم. این اسامی میتواند برای بهتر و کامل کردن دادگان استفاده شده و سبب شود که به ازای اسم فیلم ما نقد نوشته شده را نگه داری کنیم.

```
import re
movie_names=[]
for i in movie_nouns:
    for j,item in enumerate(i):
        # Use regular expression to find the movie name before the word "review"
        match = re.search(r"^(.*?)\s*review", item, re.IGNORECASE)
        if match:
            if len(i[:j])>0:
                temp=((i[:j].pop().strip())+" " + match.group(1).strip())
            else:
                temp= match.group(1).strip()
            movie_names.append([temp])
np.array(movie_names)
```

	0
0	godzilla x kong the new empire
1	mary poppins
2	kung fu panda 4
3	ryuichi sakamoto opus
4	the persian version
5	the devil