



دانشکده مهندسی کامپیوتر

تمرین درس

پردازش زبان طبیعی

استاد درس: دکتر مینایی

سجاد رضانی 96471298

یاسمین مدنی 97532265

نیم سال دوم

سال تحصیلی ۰۱-۰۰

فاز ۱

۱.۱ مقدمات

در این پروژه ما قصد داریم بتوانیم اخبار را بر اساس عنوان آنها دسته بندی کنیم موضوع داده تیتتر خبرهای موجود در سایت های خبری و دسته بندی آن است. که ما به دو زبان فارسی و انگلیسی دیتا جمع آوری کرده ایم از این جهت که اگر بخواهیم مدل خوبی را پیشنهاد دهیم باید به گونه ای باشد که قابلیت تعمیم پذیری به چند زبان را داشته باشد از آنجا که داده ها به صورت تیتتر خبر هستند پس زبان آنها رسمی بوده و محاوره ای نیست.

۲.۱ داده

داده های گردآوری شده در دو زبان فارسی و انگلیسی و همینطور به دو صورت دیتا های اسکریپ شده از وب سایت های خبری و داده های استاندارد برای این تسک بوده. در قسمت داده های اسکریپ شده از ابزار Scrapy برای جمع آوری داده ها استفاده شده و چهاراسپایدر برای وب سایت های متفاوت نوشته شده است. که به صورت خلاصه در زیر بیان شده است

- HuffingPostSpider برای استخراج داده های سایت [huffpost](#) که شامل اخبار از دسته های گوناگون می باشد.
- Techcrunch Scrapper: اسپایدر برای استخراج داده های سایت تک کرانچ که شامل اخبار مرتبط با تکنولوژی است
- Nytimes Rss Spider: اسپایدر پارس کننده اخبار سایت New York times

- اسپایدر سایت خبر آنلاین: برای یکی از مهم ترین وب سایت های خبری ایران هم یک اسپایدر نوشته ایم

جدا از موارد بیان شده از دیتا ست های آماده برای این تسک نیز استفاده شده است که به شرح زیر است

- News Category Dataset که دیتاستی از همان سایت huffpost که بخشی از دیتا را از آن به دست آورده ایم می باشد برای گرفتن دیتاست می توانید به [اینجا](#) مراجعه فرمایید
- BBC news Data: داده های اخبار سایت BBC برای گرفتن دیتاست می توانید به [اینجا](#) مراجعه فرمایید
- دیتا ست دیگر داده های اخبار اقتصادی است که در کگل منتشر شده است که می توان در [اینجا](#) مشاهده کنید

۱.۲.۱ جمع آوری دیتا

در این قسمت توضیح کوتاهی درمورد نحوه جمع آوری دیتا و اسپایدر های نوشته شده می پردازیم

همانطور که گفته شد با استفاده از فریم ورک Scrapy به راحتی می توان اسپایدر برای وب سایت های متفاوت نوشت و سپس با نوشتن پایپلاین مناسب این دیتا ها را در فرمت مناسب نگه داری کرد. که در اینجا یک نمونه از اسپایدر ها را بررسی می کنیم با استفاده از این فریم ورک ابتدا یک کلاس می نویسیم که با ارث بری از کلاس های موجود تنها کافی است چند متد آن را پیاده سازی کنیم برای این کار به مثال گردآوری داده از وب سایت خبر آنلاین توجه کنید

ابتدا لینک و فرمت پارامتر های ان در متغیر urls آورده شده است که در ان دو پارامتر tp که در واقع ایدی موضوع است و pi که شماره صفحه است را مشاهده می کند که برای شروع ۱۰۰۰ تا از ایدی ها برای موضوعات متفاوت داده شده است. سپس هر صفحه ایی که درخواست داده می شود پس از گرفتن html ان به تابع parse داده می شود تا قسمت های مورد نیاز از آن استخراج شود که با نگاه کردن به فرمت صفحه و استخراج css selector های مناسب این موارد استخراج می شود و تا صفحه مشخص شده در page limit از اخبار این دسته بندی آورده می شود.

و در اخر هم یک ابجکت از نوع Title Item برگردانده می شود

برای نوشتن این موارد هم یک pipeline ارایه شده است که در فرمت جیسون دیتا ها را بنویسید که کد آن هم در زیر ارایه شده است. برای مشاهده دقیق تر این موارد به [ریپو](#) مراجعه شود

```

1 class KhabarOnlineScraper(scrapy.Spider):
2     name = "khabaronline"
3     page_limit = 200
4     def start_requests(self):
5         urls = [
6             f"https://www.khabaronline.ir/archive?pi=1&tp={i}"
7             for i in range(10000)
8         ]
9         for url in urls:
10
11             yield scrapy.Request(url=url, callback=self.parse)
12     def parse(self, response):
13         for item in response.css('section[id="box202"]'):
14             title = item.css("h3").css("a::text")[0].get()
15             category = item.css("p").css("span").css("a::text").get()
16             yield TitleItem(title=title, category=category)
17         try:
18             if int(re.search("pi=(\d+)", response.url)
19                 .groups(1)) < self.page_limit:
20                 yield scrapy.Request(
21                     re.sub(
22                         "pi=(\d+)",
23                         lambda exp: "pi={}".format(int(exp.groups()[0]) + 1),
24                         response.url,
25                     ),
26                     self.parse,
27                 )
28         except:
29             pass

```

```

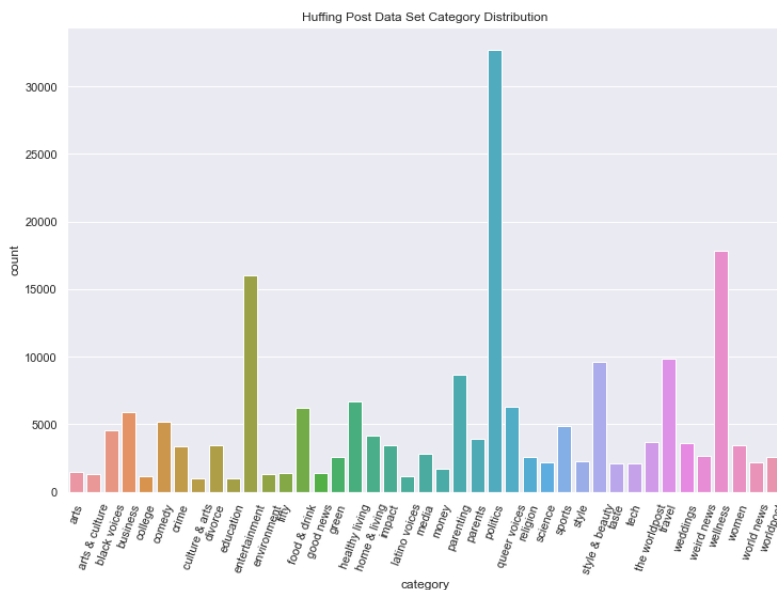
1 class JsonWriterPipeline:
2
3     def open_spider(self, spider):
4         self.file = open(ITEM_OUTPUT_PATH, 'w', encoding='utf-8')
5
6     def close_spider(self, spider):
7         self.file.close()
8
9     def process_item(self, item, spider):
10         line = json.dumps(ItemAdapter(item).asdict(), ensure_ascii=False) + '\n'
11         self.file.write(line)
12         return item

```

۲.۲.۱ دیتاست های آماده

news-category-dataset

این مجموعه داده شامل حدود ۲۰۰ هزار عنوان خبری از سال ۲۰۱۲ تا ۲۰۱۸ است که از HuffPost به دست آمده است. هر عنوان خبری یک دسته بندی مربوطه دارد که دسته بندی ها را در شکل زیر مشاهده می کنید.



Dataset: BBC

شامل ۲۲۲۵ داده از وبسایت خبری بی بی سی در پنج حوزه موضوعی از سال ۲۰۰۴ تا ۲۰۰۵ است. دارای ۵ کلاس از نوع (کسب و کار، سرگرمی، سیاست، ورزش، فناوری) می باشد

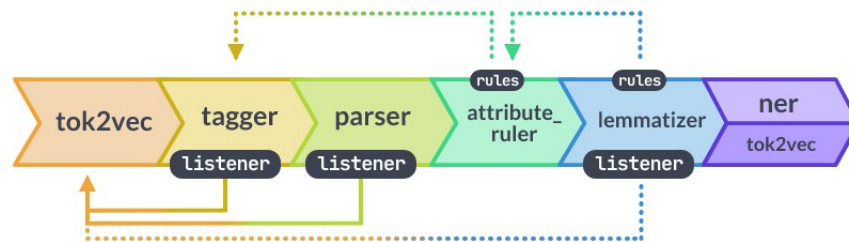
Category		
Heading	Article	
Business	510	510
Entertainment	386	386
Politics	417	417
Sport	511	511
Tech	401	401

US Financial News Articles

داده های مربوط به اخبار اقتصادی جمع آوری شده از سایت های Bloomberg.com, CNBC.com, reuters.com, wsj.com, fortune.com

۳.۱ پیش پردازش

در این مرحله نیاز داریم تا داده های جمع آوری شده را تمیز و پیش پردازش کنیم که روش های گوناگون و تسک های مختلفی از جمله ریشه گیری کلمات، حذف استاپ وردها و ... دارد. در اینجا از پایپ لاین آموزش داده شده ی `en_core_web_sm` استفاده می کنیم این پایپ لاین ابتدا متن را نشانه گذاری می کند تا یک شی `Doc` تولید کند. سپس `Doc` در چندین مرحله مختلف پردازش می شود. خط لوله معمولاً شامل یک `tagger` یک `lemmatizer`، یک `parser` و یک `entity recognizer` است. که هر یک از این بخش ها متن پردازش شده را به بخش بعدی می فرستد. برای مشاهده داک مدل به [اینجا](#) مراجعه کنید.



در ادامه در تابع زیر باقی عملیات لازم جهت پاکسازی متن را انجام می دهیم. که این شامل حذف استاپ وردها، اعداد، فاصله ها... می شود.

```

1 def clean_doc(d):
2     doc = []
3     for t in d:
4         if not any([t.is_stop, t.is_digit,
5                     not t.is_alpha, t.is_punct, t.is_space,
6                     t.lemma_ == '-PRON-']):
7             doc.append(t.lemma_)
8     return ' '.join(doc)
9
10
11 def preprocess(articles):
12     iter_articles = (article for article in articles)
13     clean_articles = []
14     for i, doc in
15         enumerate(nlp.pipe(iter_articles, batch_size=100, n_process=8), 1):
16         if i % 1000 == 0:
17             print(f"{i / len(articles):2%}.", end=" ", flush=True)
18             clean_articles.append(clean_doc(doc))
19     return clean_articles

```

برای داشتن دید بهتر تعدادی مثال از تیتراخبار پیش و پس از پیش پردازش در تصاویر زیر به نمایش گذاشته شده است.

داده پیش و پس از پردازش

raw: Captive Medic's Bodycam Shows Firsthand Horror Of Mariupol

cleaned: Captive Medic Bodycam show Firsthand Horror Mariupol

داده پیش و پس از پردازش

raw: Russia Is Firing Its Senior Commanders. What Does That Mean For Ukraine War?

cleaned: Russia fire senior commander mean Ukraine War

داده پیش و پس از پردازش

raw: LinkedIn Settles With U.S. Over Alleged Pay Discrimination

cleaned: LinkedIn Settles Alleged Pay discrimination

داده پیش و پس از پردازش

raw: Starbucks Workers Have Unionized More Than 50 Stores In The U.S.

cleaned: Starbucks Workers unionize Stores

یک مثال از شرایطی که می تواند در تصمیم گیری مدل در آینده تاثیر بگذارد

raw: LinkedIn Settles With U.S. Over Alleged Pay Discrimination

cleaned: LinkedIn Settles Alleged Pay discrimination

۴.۱ آمار مربوط به دیتاهای جمع آوری شده

از آنجا که دیتاهای مختلفی جمع آوری کرده ایم در بخش زیر تعدادی نمودار مربوط به هر یک از این دیتاستها را به نمایش گذاشته ایم اما از آنجا که انتظار می رود داده های جمع شده از سایت huffpost و دیتاست آن پاسخ بهتری در فازهای بعدی برای ما فراهم کند آمار خواسته شده در داک پروژه را در مورد این دیتاست نوشته ایم. برای بدست آوردن این اماره ها از توابع موجود در nltk استفاده شده است.

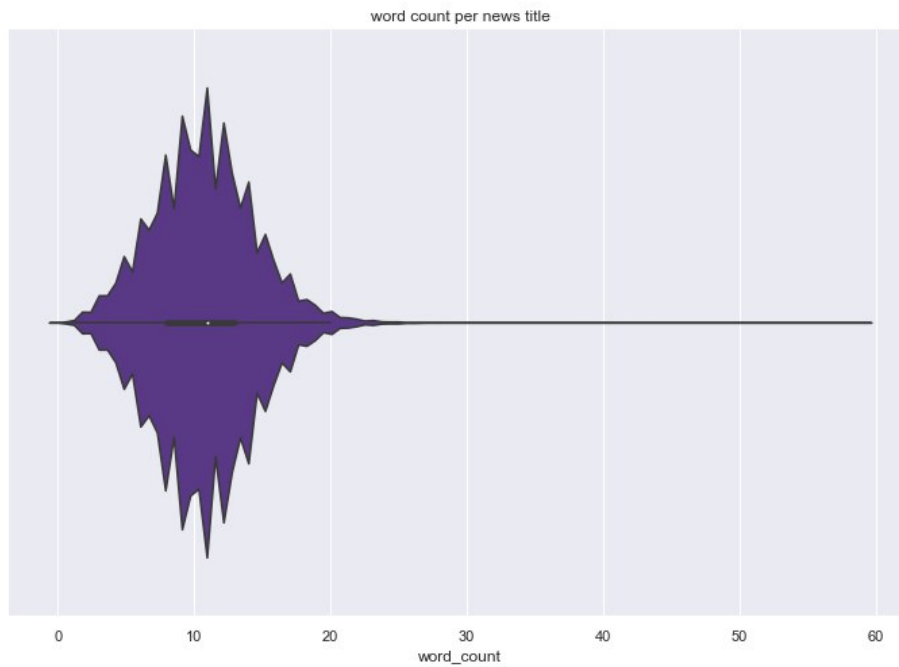
برای تفکیک جملات از `sent_tokenize` استفاده شده است. که در صورت پیش فرض با استفاده از `punctuation` ها می تواند جملات را تفکیک کند برای مطالعه بیشتر به [اینجا](#) مراجعه کنید

برای تفکیک کلمات از `word_tokenize` استفاده شده است. که در صورت پیش فرض با استفاده از `TreebankWordTokenizer` استفاده می کند.

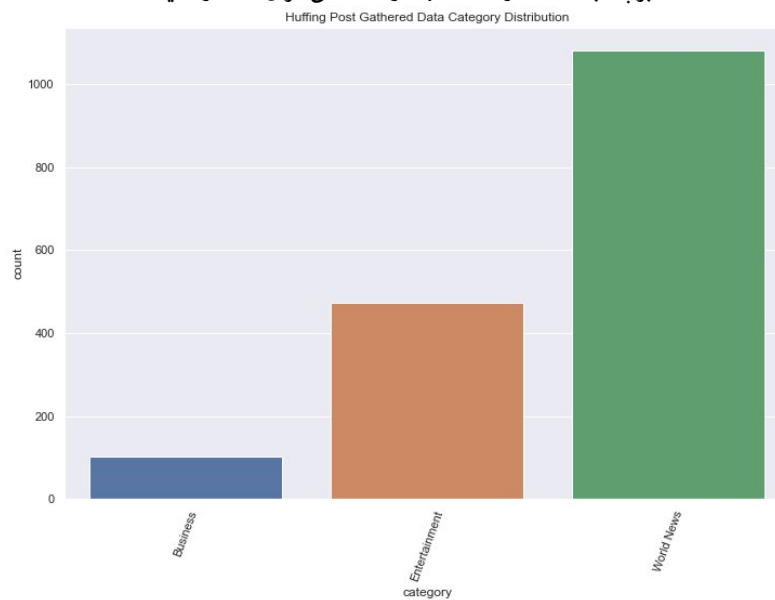
۱.۴.۱ huffpost

- تعداد واحد داده : ۲۰۱۰۵۹ عنوان خبر
- تعداد جملات : ۲۰۹۵۵۷
- تعداد کلمات ۲۱۴۶۸۹۱
- تعداد کلمات منحصر به فرد: ۶۴۷۳۲ (بعد تمیز کردن)
- تعداد کلمات منحصر به فرد: ۷۴۵۰۵ (قبل تمیز کردن)

تعداد توکنهای موجود برحسب هدلاین



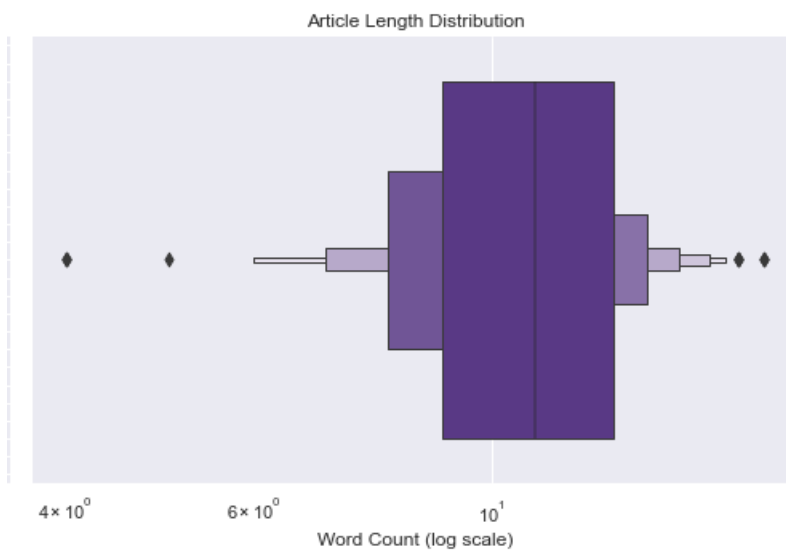
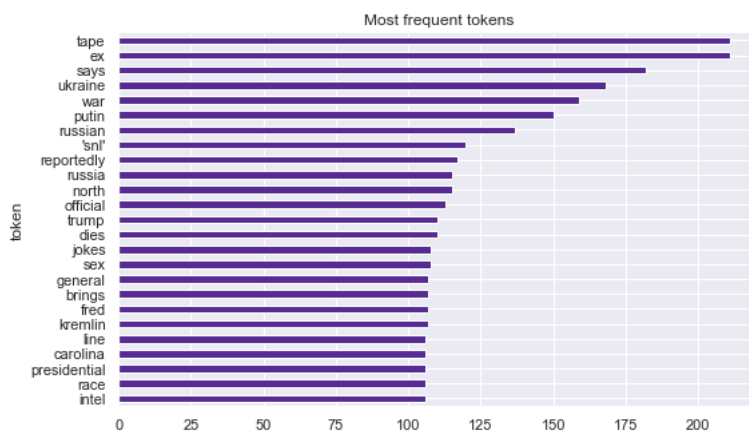
برچسب داده ها و تعداد آنها در داده های کراال شده از سایت



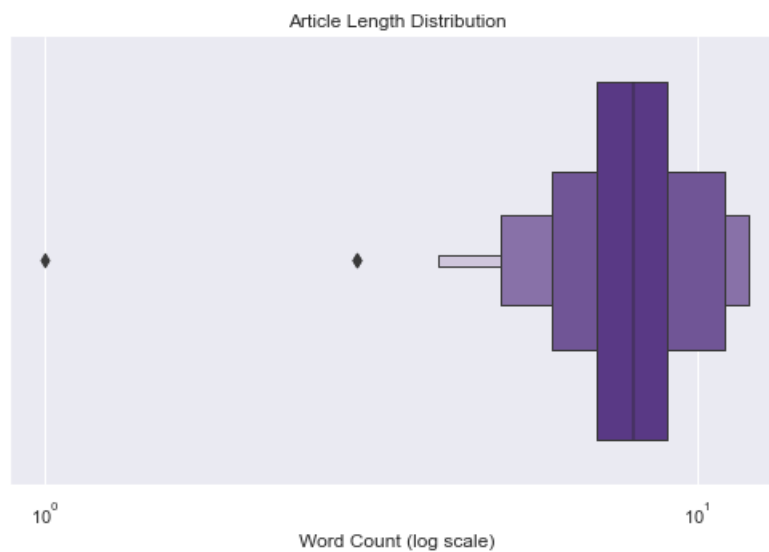
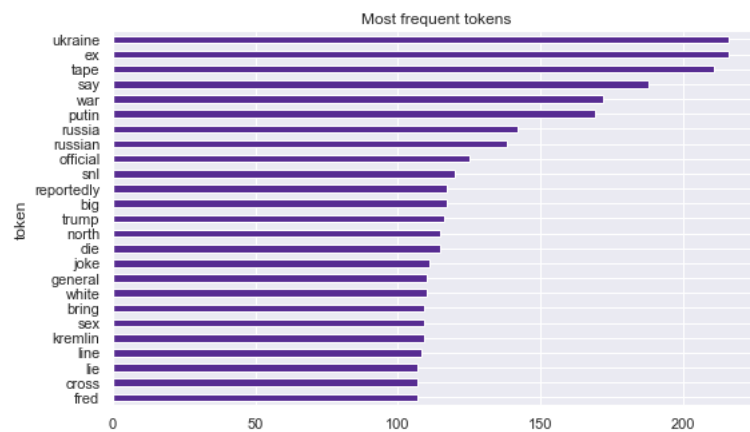
۴.۱. آمار مربوط به دیتاهای جمع آوری شده

۱۰

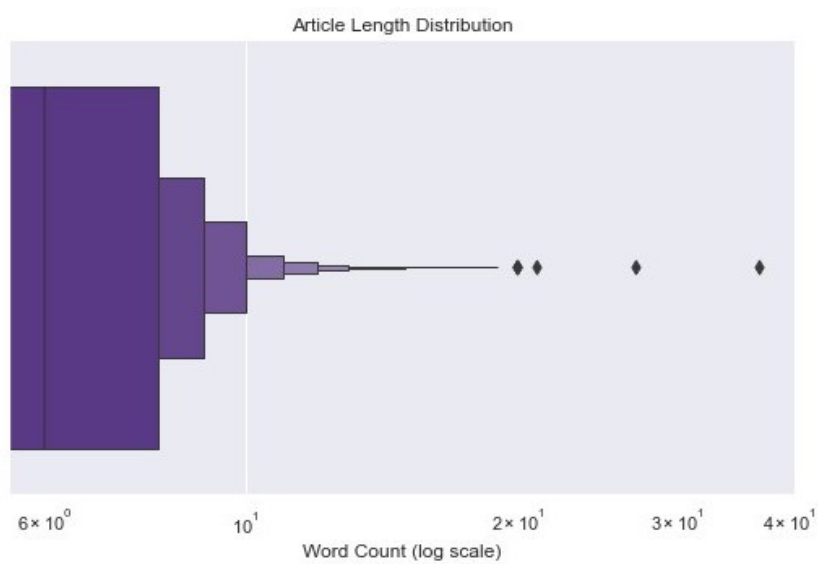
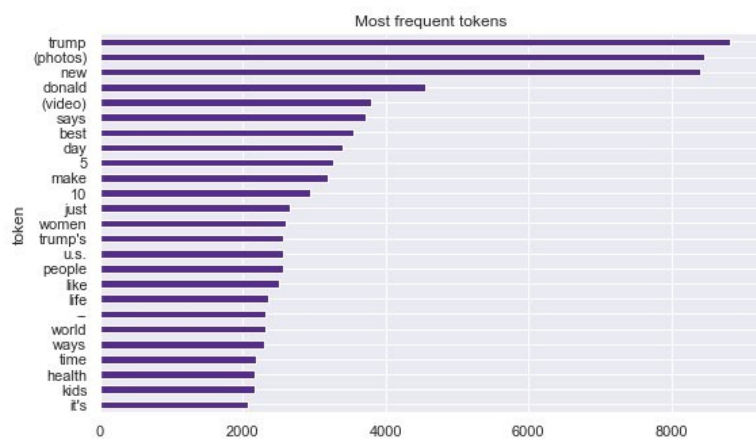
در شکل زیر نمودارهای توکن های پرتکرار و طول جملات داده های جمع آوری شده را مشاهده می کنیم.

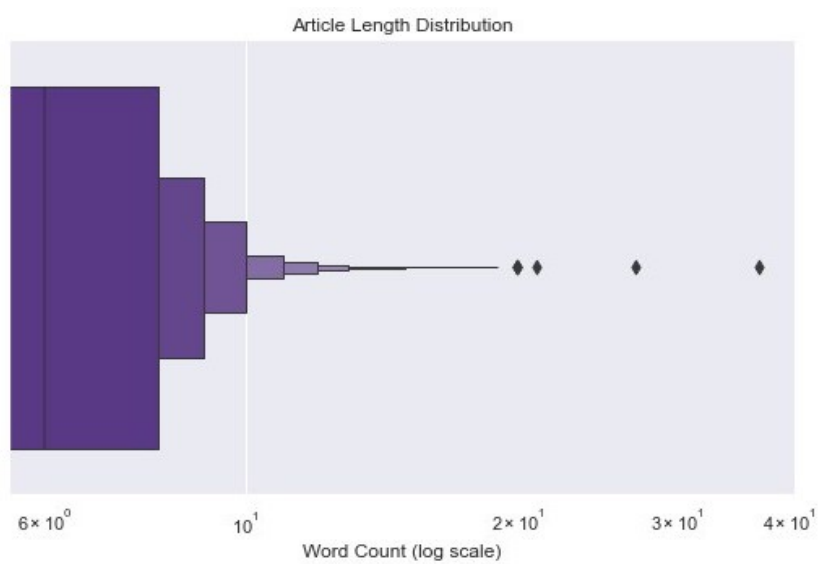
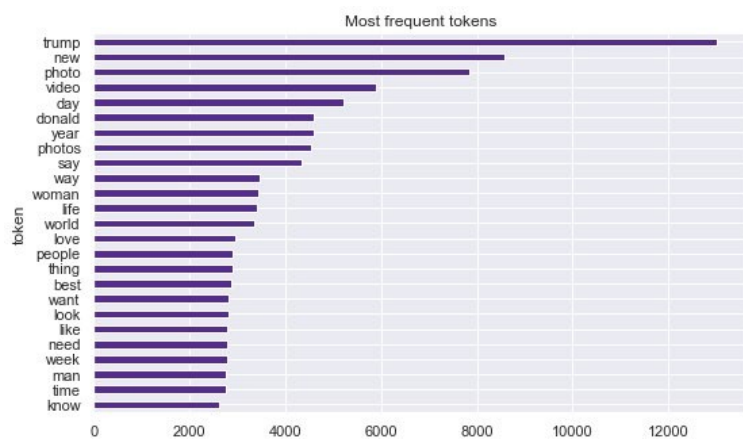


همچنین در شکل زیر نمودارهای مربوطه را پس از عملیات پیش پردازش می توان مشاهده کرد



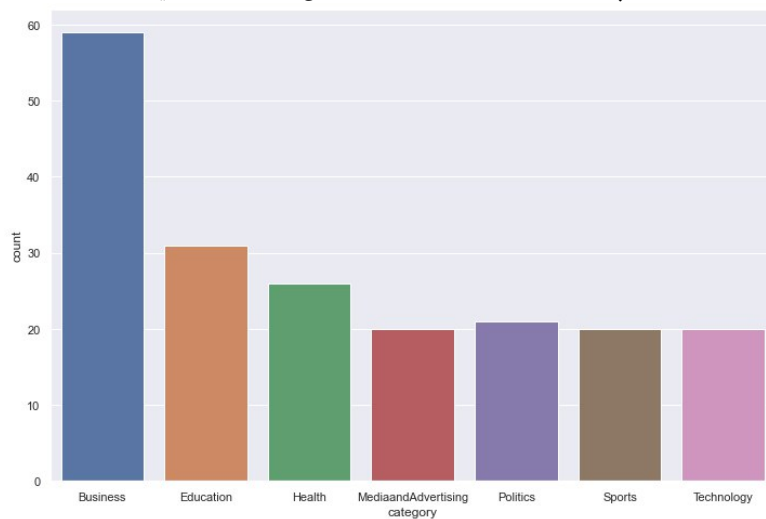
به علاوه بخشی از داده های ما از طریق دیتاست توضیح داده شده پیش از این تامین شده از این رو این عملیات پیش پردازش را روی داده های موجود در دیتاست نیز اجرا می کنیم همان طور که در نمودارهای زیر مشاهده می شود داده های موجود در دیتا بیس حاوی مقادیر اطلاعات از قبیل اعداد و ... اند که با انجام عملیات پیش پردازش حذف شده و اولویت در ترتیب توکن های پر تکرار تغییر می کند.





۲.۴.۱ nyt

برچسب داده ها و تعداد آنها در داده های کرال شده از سایت



۳.۴.۱ خبرآنلاین

برچسب داده ها و تعداد آنها در داده های کرال شده از سایت

