



به نام خدا

قسمت تحلیل عواطف و احساسات و هشتگ گذاری
پروژه نهایی درس تحلیل و سیستم‌های داده‌ای حجیم

استاد درس:

دکتر نعمت بخش

بهنام صوفی

۴۰۱۳۶۱۴۰۲۸

تیر ۱۴۰۲

لینک گیت‌هاب:

در مسیر زیر :

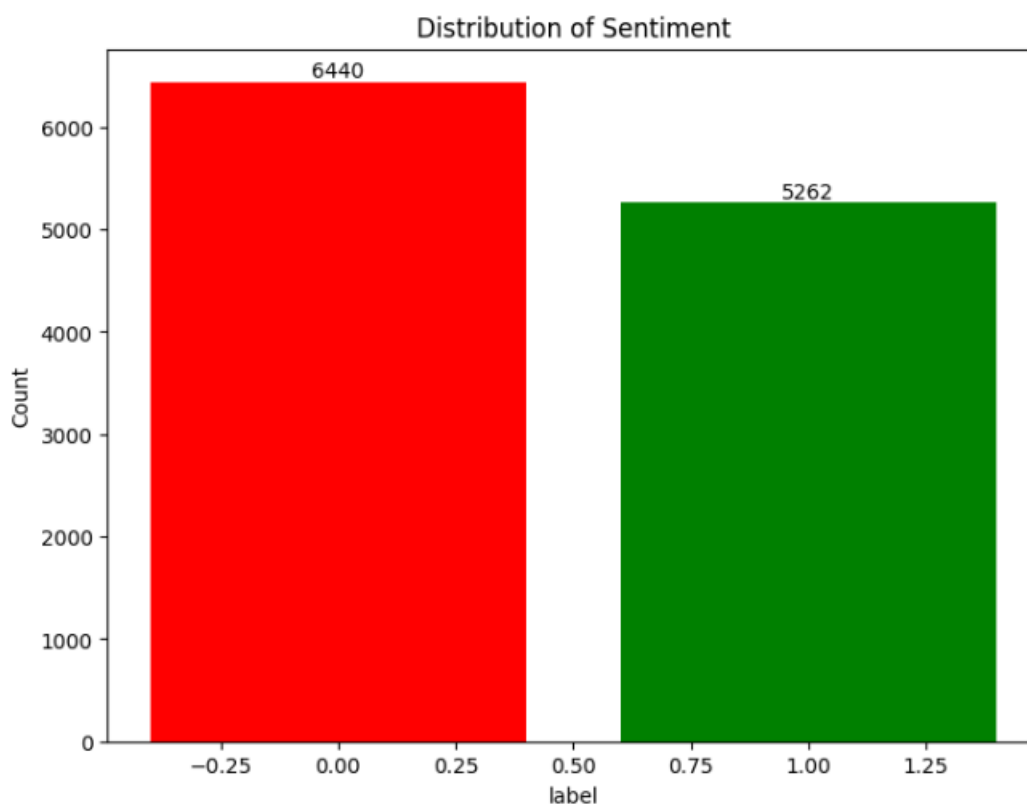
<https://github.com/bscoder9595/BigDataProjects/>

فایل bigdata_finalproject.rar

کار ایجاد مدل تحلیل احساسات و ایجاد برچسب در ۶ مرحله انجام شد، در ادامه مراحل را توضیح می‌دهیم:

۱. گرفتن داده‌ها

ما داده‌ها را از فایل CSV می‌گیریم و در قالب دیتافریم اسپارک ذخیره می‌کنیم، در داده‌هایی که ما داریم ۱ به عنوان خبر مثبت و ۰ به معنی خبر منفی است. رنگ سبز نشان‌دهنده اخبار مثبت و رنگ قرمز نشان‌دهنده اخبار منفی است. در شکل زیر توزیع داده‌ها را مشاهده می‌کنید:



۲. انجام اعمال پیش پردازش

چهار عمل اصلی که در پیش پردازش انجام می شوند عبارتند از :

- توکن سازی

در این مرحله، متن ورودی (اطلاعات متنی مانند یک مقاله یا پاراگراف) به صورت جملات یا واحدهای کوچکتر تقسیم می شود، که به عنوان "توکن ها" شناخته می شوند. توکن ها ممکن است کلمات، نقاط، ویرگول ها، علائم نگارشی و ... باشند. این کار اولیه مرحله ای اساسی است که متن را به یک سری از توکن ها تبدیل می کند.

- حذف Stop Words

کلمات توقف (Stop Words) کلمات معمولاً رایج در زبان هستند و اغلب اطلاعات کمی را به ما می دهند. این کلمات شامل حروف اضافه، افعال متعلق به بودن، حروف ربط و ... می شوند. البته ما علائم نشانه گذاری را هم در این فایل قرار دادیم. در این مرحله، توکن هایی که کلمات توقف هستند، از متن حذف می شوند. این کار به کاهش تعداد واژگان غیر ضروری کمک می کند و می تواند در بهبود عملکرد مدل های یادگیری ماشینی کمک کند.

- تبدیل توکن ها به وکتورها

در این مرحله، توکن های پاک سازی شده به بردارهای عددی تبدیل می شوند. برای انجام این کار، از روش های مختلفی مانند CountVectorizer و TF-IDF استفاده می شود. در اینجا، از روش CountVectorizer استفاده شده است. این روش تعداد تکرارهای هر توکن را در متن مشخص می کند و این تعداد تکرار را به عنوان یک بردار عددی نمایش می دهد.

- محاسبه TF-IDF

TF-IDF یک روش محاسبه ویژگی‌های مهم در متن‌ها است که با استفاده از دو مفهوم "تعداد تکرار کلمه در سند (TF)" و "تعداد سندهایی که کلمه را شامل می‌شوند (IDF)" عمل می‌کند. TF-IDF برای اندازه‌گیری اهمیت یک کلمه در یک متن نسبت به مجموعه اسناد کلی استفاده می‌شود. با محاسبه این مقادیر، بردارهای ویژگی TF-IDF برای هر سند (متن) ایجاد می‌شود که برای استفاده در مدل‌های یادگیری ماشینی مناسب هستند. این ویژگی‌ها معمولاً به عنوان ورودی به مدل‌های یادگیری ماشینی تحویل می‌شوند تا بتوانند روی متن‌ها مسائلی مانند دسته‌بندی، خوشه‌بندی و ... را انجام دهند.

در داده‌های موجود، ما فقط به ستون label، article و title نیاز داریم. در قدم اول باید عمل توکن‌سازی را بر روی ستون article انجام دهیم و کلمات موجود در آن متن را به صورت جدا بنویسیم، در قدم بعدی کلماتی که معروف به stop هستند و در یک فایل جدا آنها را آماده کردیم را از متن حذف کنیم. در مرحله بعدی کلمات جدا شده را وکتور می‌کنیم و در آخر TF-IDF را برای آن‌ها محاسبه می‌کنیم.

۳. فراخوانی ۵ نوع الگوریتم و انتخاب یکی از آن‌ها

از ۴ الگوریتم Logistic Regression، Svm، Random Forest و NaiveBayes و شبکه cnn برای ایجاد مدل تحلیل احساسات استفاده شد که نتایج زیر را داشت و با توجه به اینکه cnn در بستر pyspark انجام نشده است، svm انتخاب شده است.

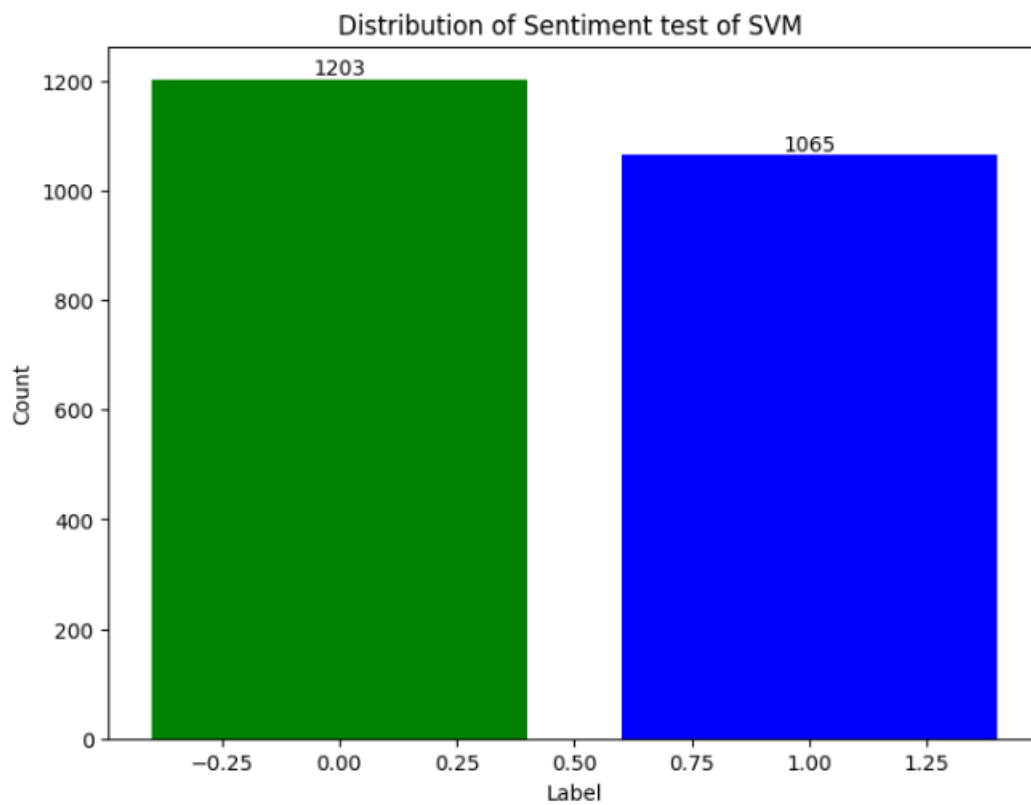
ایجاد مدل تحلیل احساسات

ما از همه الگوریتم‌ها استفاده کرده در جدول زیر دقت هر کدام را مشاهده می‌کنید:

Accuracy	الگوریتم‌ها
0.743	Svm
0.701	Logistic Regression
0.663	Random Forest
0.717	NavieBayes
0.973	Cnn

نمودار زیر توزیع پیش‌بینی SVM روی داده‌های تست را نشان می‌دهد:

رنگ آبی نشان‌دهنده تعداد داده‌هایی که مثبت در نظر گرفته است و رنگ سبز تعداد نظراتی که منفی در نظر گرفته است.



۴. ذخیره مدل و لود دوباره آن برای تست

به دلیل اینکه وظایف تقسیم شده، مدل ایجاد شده را ذخیره کرده و دوباره لود کرده و تست می‌کنیم که مطمئن شویم مدل به درستی کار می‌کند. فراخوانی مدل ذخیره‌شده و تست دوباره مدل را در تصاویر زیر مشاهده می‌کنید:

```
newspd = test_news_df.toPandas()
newspd
```

	article
0	... در اطلاعیه‌های منتشر شده امروز در سامانه گدال
1	... در اطلاعیه‌های منتشر شده امروز در سامانه گدال
2	... همچنین در اطلاعیه‌های امروز سامانه گدال شرکت س

```
lr_loaded = LogisticRegressionModel.load('/content/drive/MyDrive/csvFiles/Logistic_model')

rf_loaded= RandomForestClassificationModel.load('/content/drive/MyDrive/csvFiles/random_model')

svm_loaded = LinearSVCModel.load('/content/drive/MyDrive/csvFiles/svm_model')
nv_loaded = NaiveBayesModel.load('/content/drive/MyDrive/csvFiles/NaiveBayes_model')

lr_loaded_predictions = lr_loaded.transform(test_news_df)
rf_loaded_predictions = rf_loaded.transform(test_news_df)
svm_loaded_predictions = svm_loaded.transform(test_news_df)
nv_loaded_predictions= nv_loaded.transform(test_news_df)
```

```
lr_loaded_sentiments = ["Positive" if label == 1 else "Negative" for label in lr_loaded_predictions.select("predictor")]
rf_loaded_sentiments = ["Positive" if label == 1 else "Negative" for label in rf_loaded_predictions.select("predictor")]
svm_loaded_sentiments = ["Positive" if label == 1 else "Negative" for label in svm_loaded_predictions.select("predictor")]
nv_loaded_sentiments = ["Positive" if label == 1 else "Negative" for label in nv_loaded_predictions.select("predictor")]

print("Logistic Regression Sentiments_loaded:", lr_loaded_sentiments)
print("Random Forest Sentiments_loaded:", rf_loaded_sentiments)
print("SVM Sentiments_loaded:", svm_loaded_sentiments)
print("NV Sentiments_loaded:", nv_loaded_sentiments)

Logistic Regression Sentiments_loaded: ['Negative', 'Negative', 'Negative']
Random Forest Sentiments_loaded: ['Negative', 'Negative', 'Negative']
SVM Sentiments_loaded: ['Negative', 'Negative', 'Negative']
NV Sentiments_loaded: ['Negative', 'Negative', 'Negative']
```

۵. ایجاد برچسب بر روی داده‌ها

ما به دو روش برچسب گذاری را انجام داده ایم روشی که اینجا در مورد آن توضیح خواهیم داد روش دوم است. تابعی داریم تحت عنوان `extract_keywords` که وظیفه دارد که از لیست کلمات دریافتی، کلمات کلیدی را استخراج کند. این تابع ابتدا لیست کلمات را به عنوان ورودی دریافت می‌کند، سپس با استفاده از لیست `stopwords`، کلمات ناخواسته را حذف می‌کند. سپس تعداد تکرار هر کلمه را محاسبه کرده و کلمات را بر اساس تعداد تکرار به صورت نزولی مرتب

می‌کند. سپس لیست stopwords_list2 نیز به عنوان لیست دیگری از کلمات ناخواسته برای حذف اعمال می‌شود. در نهایت، تا سه کلمه کلیدی با تعداد تکرار بالاتر به عنوان نتیجه برگشت داده می‌شوند.

۶. ذخیره داده‌ها در قالب دیتافریم و فایل CSV

داده‌ها در نهایت بصورت دیتافریمی از عنوان شرکت، خبر، پیش‌بینی این‌که مثبت یا منفی است و دو مجموعه هشتگ ذخیره می‌شود. دیتافریم زیر نشان‌دهنده دیتافریم خروجی است:

hashtag2	hashtag1	prediction	article	title
نرخ باشد, مذکور	میباشد, نرخ مذکور, مبلغ میلنگ	0.0	...به گزارش کدال نگر بورس24، شرکت صنایع ریخته گ	«افزایش نرخ های فروش »خریفت
ماهه, فروش ,	ماهه, فروش, مبلغ, بورس	1.0	...به گزارش کدال نگر بورس24، شرکت صنعتی پارس خ	روند فروش "لنزر" را اینجا ببینید
ماهه, فروش ,	ماهه, فروش, مبلغ, بورس	1.0	...به گزارش کدال نگر بورس24، شرکت شیر یاستوریزه	"بررسی روند فروش "فشادر
عرضه کارگزاری, توسعه	عرضه بورس, کارگزاری, توسعه, سپند	0.0	...به گزارش کدال نگر بورس24، شرکت کارگزاری توسعه	بلوک 7 درصد «گنگین» عرضه می شود
خرداد, کرد ,	خرداد, ارزش, بورس	1.0	...به گزارش کدال نگر بورس24، شرکت سرمایه گذاری ا	تحقق سود 6 میلیارد تومانی «سنوین» در خرداد ماه
...
سود, مجتمع ,	سود, بورس, مجتمع, معادن, نکدار	1.0	...به گزارش کدال نگر بورس24، شرکت مجتمع معادن مس	بررسی عملکرد «نکدار» تا پایان شهریور
سود, صنعتی, شیمیایی	سود, بورس, صنعتی, شیمیایی, رنگین	1.0	...به گزارش کدال نگر بورس24، شرکت صنعتی و شیمیایی	«تحقق سود 18 میلیارد تومانی در «شرنگی
شهریور, ارزش, سرمایه	شهریور, ارزش, بورس, سرمایه, توسعه	0.0	...به گزارش کدال نگر بورس24، شرکت سرمایه گذاری ت	آیا از تغییرات پرنفوذی «ومعادن» خبر دارید؟
شرکت , صنایع, سیمان	بورس, صنایع, سیمان, غرب, سود	1.0	...به گزارش کدال نگر بورس24، شرکت صنایع سیمان غ	طت رشد سود «سغرب» چه بود؟
مهر, گروه, توسعه	مهر, بورس, گروه, توسعه, مالی	1.0	...به گزارش کدال نگر بورس24، شرکت گروه توسعه مال	تحقق سود 19 میلیارد تومانی «ومهان» در مهر ماه

2269 rows x 5 columns

در نهایت دیتافریم را در یک فایل CSV برای استفاده در مرحله بعدی، ذخیره می‌کنیم.