

به نام خدا

گزارش تمرین دوم درس NLP

دانشجو: محمد قربانی-۹۷۱۳۱۰۹۹

دانشگاه صنعتی امیرکبیر - دانشکده مهندسی کامپیوتر

بهار ۹۷-۹۸

روش اول (word2vec و استفاده از میانگین ساده):

در این روش ابتدا یک پیش پردازش اولیه بر روی داده ترین انجام صورت گرفت و کلمات ایست که تنها سربار محاسباتی به سیستم اضافه می‌کند حذف شدند. بعد از حذف کلمات ایست، مطالب را بر روی یک فایل ریخته و سپس بر روی آنها یک word2vec ساختیم. در ساخت مدل مقدار min_count را برابر با ۱ در نظر گرفتیم که مشکل کلمات دیده نشده بوجود نیاید. سپس با استفاده از این مدل بردار کلمات هر داکيومنت را به دست آوردیم و سپس این بردارها را برای هر داکيومنت جمع کردیم و تقسیم بر تعداد کلمات آن داکيومنت‌ها کردیم و این بردار را به عنوان بازنمایش آن داکيومنت در نظر گرفتیم. بعد از ایجاد بردار بازنمایی هر داکيومن این بردارها را به الگوریتم کلاسترینگ خود دادیم تا مدل کلاسترینگ خود را ایجاد کنیم. بعد از ایجاد خوشه‌بند، بردار هر داکيومنت داده ترین خود را به خوشه‌بند دادیم تا آن را خوشه‌بندی کند. سپس بیشترین تعداد موضوع در هر خوشه را به عنوان لیبل آن خوشه در نظر گرفتیم و سپس به محاسبه NMI، Accuracy، F-Measure و V-Measure پرداختیم.

| داده تست | داده ترین | معیار |
|----------|-----------|-----------|
| 0.58 | 0.56 | NMI |
| 0.72 | 0.72 | Accuracy |
| 0.74 | 0.73 | F-Measure |
| 0.58 | 0.56 | V-Measure |

NMI یکی دیگر از معیارهایی است که برای سنجش کیفیت خوشه‌بندی استفاده می‌شود. این معیار با در نظر گرفتن خوشه‌بندی واقعی آیتم‌ها و همچنین خوشه‌بندی تولید شده توسط الگوریتم شباهت بین این دو مجموعه از خوشه‌ها را بررسی می‌کند. این معیار بررسی می‌کند که این دو دسته از خوشه‌ها تا چه اندازه به هم مشابه و تا چه اندازه نسبت به هم متفاوت هستند. این معیار مقداری بین ۰ و ۱ دارد که مقدار ۱ بیانگر خوشه‌بندی ایده آل می‌باشد. مزیت معیار NMI این است که نسبت به تعداد خوشه‌های در نظر گرفته شده نرمال می‌شود.

نتایج حاصل از این روش نشان می‌دهند که این روش دقت نسبتاً خوبی داشته است.

روش دوم (word2vec و استفاده از میانگین وزن دار):

در این روش بردار هر کلمه در وزن tfidf آن کلمه ضرب شده است. مقدار tfidf به صورت زیر محاسبه می‌شود:

$$tf_{idf} = tf * idf$$

$$tf = \begin{cases} 1 + \log count(t, d) & \text{if } count(t, d) > 0 \\ 0 & \text{else} \end{cases}$$

$$idf = \log\left(\frac{N}{df_i}\right)$$

که در آن N بیانگر تعداد کل اسناد موجود در پیکره است و df_i بیانگر تعداد داکيومنت‌هایی که کلمه مورد نظر را دارند است. استفاده از idf موجب می‌شود که مشکل کلمات ایست که سربار محاسباتی دارند و هیچ گونه اطلاعات مفیدی نمی‌دهند حل شود. کلمات ایست حاصل idf کمی می‌گیرند و این باعث می‌شود که در محاسبات ارزش کمی بگیرند.

نتایج حاصل از این روش به صورت زیر می‌باشد. در مقایسه با روش اول کمی کاهش اندازه معیارها را نشان می‌دهد. دلیل این امر این است که ما در روش اول خود در ابتدا در یک پردازش اولیه کلمات ایست را از پیکره حذف نمودیم و سپس به آموزش مدل‌های خود پرداختیم اما در این روش این کار به اثرات حاصل از فرمول idf سپردیم که نتایج نشان می‌دهد که بهتر است حذف کلمات ایست را به عنوان بخشی از اقدامات پیش پردازشی خود انجام دهیم.

| داده تست | داده ترین | معیار |
|----------|-----------|-----------|
| 0.54 | 0.52 | NMI |
| 0.67 | 0.66 | Accuracy |
| 0.66 | 0.64 | F-Measure |
| 0.54 | 0.52 | V-Measure |

روش سوم (doc2vec):

در این روش برای بازنمایی هر سند از الگوریتم doc2vec استفاده شد. این روش در واقع یک الگوریتم بدون سرپرستی می باشد که در آن برای جمله/پاراگراف/داکیومنت یک بردار ایجاد می کند. این الگوریتم در واقع الهام گرفته از word2vec می باشد که در آن برای هر کلمه یک بردار ایجاد می کند. در این روش کلمات به صورت مستقل از هم در نظر گرفته نمی شود و یک رابطه ی سیکونسی بین کلمات در نظر گرفته می شود.

| داده تست | داده ترین | معیار |
|----------|-----------|-----------|
| 0.18 | 0.17 | NMI |
| 0.50 | 0.51 | Accuracy |
| 0.47 | 0.46 | F-Measure |
| 0.18 | 0.17 | V-Measure |

این روش نسبت به دو روش قبل دارای نتایج ضعیف تری می باشد که این موضوع کمی عجیب به نظر می رسد چرا که انتظار ما این بود که در این روش نسبت به دو روش قبل نتایج بهتری مشاهده کنیم. به همین دلیل به بررسی چندین باره روند کار پرداختیم اما متوجه دلایل نتایج ضعیف نشدیم.

روش چهارم (document-term):

در این روش برای بازنمایی هر سند از ماتریس سند-کلمه استفاده شد. اندازه‌ی این ماتریس برابر با مقداری در حدود 8600×65000 شد. برای بازنمایی از کل سندهای آموزش و ارزیابی استفاده شد. بعد از این که ماتریس حاصل محاسبه شد ما برای این که آن را تبدیل به یک ماتریس dense کنیم از الگوریتم svd استفاده کردیم. از آنجایی که این ماتریس دارای مقادیر صفر زیاد بود برخی از الگوریتم‌ها در زمان مناسب پاسخگو نبوده و ما مجبور به استفاده از الگوریتم‌های خاصی برای بدست آوردن ماتریس dense شدیم. الگوریتم مورد استفاده ما svds که از کتابخانه‌ی scipy است بود.

بعد از بدست آوردن ماتریس dense با استفاده از قسمت مربوط به آموزش آن، مدل خوشه‌بندی خود را آموزش دادیم. بعد از خوشه‌بندی، طبق الگوی گفته شده هر خوشه را ورچسب‌گذاری کرده و به محاسبه معیارهای خواسته شده پرداختیم. در ادامه به سراغ داده تست رفته و بازنمایی که برای هر کدام از این اسناد به دست آورده بودیم را به مدل خوشه‌بند خود داده تا ورچسب آن را پیش‌بینی کند. سپس به محاسبه معیارهای خواسته شده پرداختیم که نتایج حاصل آن به صورت زیر می‌باشد:

| داده تست | داده ترین | معیار |
|----------|-----------|-----------|
| 0.04 | 0.04 | NMI |
| 0.27 | 0.28 | Accuracy |
| 0.42 | 0.42 | F-Measure |
| 0.04 | 0.04 | V-Measure |

این روش نسبت به روش‌های قبل از نتایج ضعیف‌تری برخوردار می‌باشد که به نظر نویسنده یکی از دلایل آن می‌تواند استفاده از الگوریتم svd نامناسب باشد.