

گزارش پایانی فاز سوم پروژه

95521198

آریانا سیدپیروی

۱- با توجه به انتخاب شده های نادرست (استقلال انتخاب شده اند ولی در اصل پرسپولیس هستند).
مثلا سطر سوم فایل متنی واقع در آدرس

`path\to\git\ClsModel\NaiveBayes\group_samples\fp_esteghlal.txt`
classifier برای استقلال میتوانیم دریابیم که بعضی جمله ها کلمات کلیدی مانند پرسپولیس آمده است اما چون طول جمله زیاد شده و کلمات تعدادشان زیاد شده چند کلمه دیگر آمده اند و باعث شده اند احتمال کلمه پرسپولیس کم شود. پس این جا را می توان از نقطه ضعف های این مدل به شمار آورد. (این مورد در سطر هفتم همین فایل در مورد برانکو نیز صادق است.) با توجه به انتخاب نشده های درست (استقلال انتخاب نشده اند ولی در اصل استقلالی هستند.) همین مورد بالا در سطر بیستم فایل متنی واقع در آدرس

`path\to\git\ClsModel\NaiveBayes\group_samples\fn_esteghlal.txt`

نیز دیده می شود. همچنین در خط ۲۵ این فایل با وجود کلمه ای مثل محسن کریمی که به صورت پیوسته باید در نظر گرفته شود با اینکه در مقایسه با کلمه علی کریمی تعداد کم تری دارید ولی باید متن استقلالی باشد نه پرسپولیس ولی به دلیل احتمال بالای کریمی این جمله به پرسپولیس تشخیص داده شده. (البته بنده از فرمول `laplace smoothin` و در نظر گرفتن کلمه `<Unk>` فرمول NaiveBayes را نوشتم که در مخرج $|V| + 1$ داشت.

۲- با توجه به اینکه دیتاست من کاملا متفاوت بود من به تنهایی با تعداد کلمات توانستم نتیجه ی خوبی در این بخش کسب کنم و دقت بالای ۹۰ درصد را به دست بیاورم. به نظر من با توجه به دیتای شخصی معیار های کمی مثل تعداد کلمات آبی پوشان یا سرخ پوشان یا ... میتوانستند ویژگی های بسیار به درد بخور و حایز اهمیتی باشند که همینطور هم شد ویژگی دیگری به درد بخور به ذهنم خطور نکرد. (البته خواستم فقط از تعداد کلمات موجود در word cloud ها برای این کار استفاده کنم ولی متوجه شدم کلمات به درد نخور هم ممکن است در بی آن ها یافت شده و حذف شوند و همچنین کلمات پرارزشی که دارای تکرار کمی باشند اما وزن بالایی داشته باشند.

۳- با توجه به دقت بالا من ویژگی های دیگری اضافه نکردم.

۴- به نظر من استفاده از MaxEnt خیلی بهتر بود چون دقت محاسبات را تقریبا به بالای ۹۰٪ رساند و همچنین F1_measure تفاوت بسیار زیادی کرد که در نوع خودش عالی بود ولی با توجه به محاسبات ساده تر من خودم NaiveBayes را ترجیح می دهم چرا که سرعت بالایی دارد و تفاوت زیادی با MaxEnt ندارد ولی شاید در دیتاهای بیشتر با چالش های بیشتر مثلا جاهایی که ترتیب کلمات خیلی مهم باشند یا جمله های طولانی تری باشند استفاده از MaxEnt دقت را بسیار بالا

ببرد.البته من در فیچرهایم کلمه <Unk> را هم در نظر گرفتم چون کلمه ای است که در train data دیده نشده است.