

به نام خدا

پروژه ترم اینجانب شامل ابهام‌زدایی معنایی از 3 کلمه شیر، کرم و شانه است.

نحوه کار به این صورت بوده که ابتدا فایل‌های متنی پیکره متنی متعلق به هر یک از کلمات بر اساس تعداد جملات/پاراگراف‌ها که برچسب‌دهی شده بودند به 4 قسمت تقریباً مساوی تقسیم شده و طی 4 مرحله هربار یک قسمت از آنها مورد آزمون قرار می‌گیرد. فایل‌های متنی پیکره که نسبت به فایل‌های ارسال‌شده کمی تغییر کرده‌اند در فولدرهای shir\_4 folds ، krm\_4 folds و shane\_4 folds قرار داده شده‌اند.

فیچرهای استخراج شده از پیکره آموزش شامل فیچرهای collocational که خود شامل باهم‌آیی‌های 2تایی و 3تایی کلمات (بایگرم و ترایگرم)، کلمات قبل و بعد (بایگرم‌های) تارگت-ورد می‌شود و فیچرهای bag of words که در 3 مرحله در پنجره‌های  $\pm 5$  تا  $\pm 10$  اطراف تارگت-ورد و از  $\pm 10$  تا ابتدا و انتهای جمله مورد سنجش قرار گرفته‌اند، می‌باشد. لیست این فیچرها در فولدر total extracted features list قرار داده شده و مقادیر از فایل features\_list\_2 به داخل کد برنامه فراخوانی می‌شوند.

برای ابهام‌زدایی از هر کدام از کلمات 4 نوع کد برای آموزش و 4 کد برای آزمون نوشته شده است.

برای مثال :

برای ابهام‌زدایی کلمه کرم، کد wd\_krm\_1\_train.py که در فولدر krm\_4 folds\_codes قرار گرفته است با استفاده از اطلاعات استخراج شده از فایل متنی krm\_1\_train.txt که در فولدر krm\_4 folds قرار دارد نوشته شده و با اجرای آن یک فایل متنی به نام krm\_1\_train\_sorted\_decision\_list.txt ایجاد می‌شود که لیست تصمیم مورد استفاده در برنامه تست به نام wd\_krm\_1\_test.py را در خود جای داده است.

لیست‌های تصمیم بر اساس مقادیر بیشینه Log Likelihood ratio محاسبه شده برای هریک از فیچرها مرتب‌سازی شده‌اند.

پس از اجرای هر 4 قسمت برنامه آموزش و ایجاد 4 لیست تصمیم، 4 برنامه تست متناظر با آن‌ها را اجرا نموده و مقدار Average\_precision از میانگین 4 precision اجرا، به دست می‌آید.

نتایج حاصل از یک نمونه از اجرای 4 برنامه تست برای هر کلمه (مثلاً با نام `krm_4 tests_average_precision_exe` برای کلمه کرم) در فولدر `krm_4 folds_codes` قرار داده شده است.