

# گزارش کار تمرین عملی سری دوم درس پردازش زبان طبیعی

محمد رضا اصولی - 610395077

## مراحل انجام کار:

ابتدا به دانلود corpus های مورد نظر پرداختیم (movie-reviews و stopwords برای حذف این نوع کلمات از review ها). سپس برای به دست آوردن ویژگی های از روش bag of words استفاده کردیم، به این روش که ویژگی های هر review بر اساس کلمات موجود در آن شناخته می شود. پس ابتدا همه کلمات موجود در همه review ها را استخراج کرده، سپس stopwords و علائم نگارشی را از آن حذف کرده (این کار باعث بالا رفتن دقت می شود، به این دلیل که این کلمات در تمامی متن ها موجود بوده و اطلاعات خاصی در مورد آن متن به ما نمی دهند، پس بهتر است که حذف شوند) و در نهایت 1000 کلمه ای که بیشتر از آن ها استفاده شده را به عنوان feature استخراج کردیم. (مقدار 1000 یه می تواند متغیر باشد ولی اینکه کلمات کم تکرار از ویژگی ها حذف شوند معمولا باعث بالا رفتن دقت در این روش می شود)

سپس برای هر review ویژگی هایش را استخراج کرده و در نهایت با روش k-fold داده ها را به 5 دسته تقسیم کرده و در هر مرحله یک دسته را به عنوان تست و 4 دسته دیگر را به عنوان داده آموزشی استفاده کردیم (پس 5 دسته بند متفاوت خواهیم داشت) که برای هر دسته بند مقادیر accuracy، recall، precision و F-measure را محاسبه کردیم.

## تحلیل تاثیر پیش پردازش ها بر نتایج رده بندی:

همانطور که در بالا گفته شد حذف stopword ها و علائم نگارشی به علت نداشتن اطلاعات مضاعف باعث بالا رفتن مقدار دقت می شود.

بهترین رده بند، بعد از حذف	بهترین رده بند، قبل از حذف
os_precision: 0.846590909090909 neg_precision: 0.78125 pos_recall: 0.7525252525252525 neg_recall: 0.8663366336633663 pos_f_measure: 0.796791443850 neg_f_measure: 0.821596244131 accuracy: 0.81	pos_precision: 0.7864583333333334 neg_precision: 0.7115384615384616 pos_recall: 0.7156398104265402 neg_recall: 0.783068783068783 pos_f_measure: 0.7493796526054591 neg_f_measure: 0.7455919395465995 accuracy: 0.7475

## توضیح ویژگی ها و دلیل انتخاب آن ها:

در این روش یک ویژگی به معنی داشتن یا نداشتن یک کلمه خاص تلاقی می شود. این روش، روشی ساده اما کارآمد برای تحلیل احساسات در متون است، به این صورت که ترکیب کلمات با بارهای معنایی متفاوت می تواند تا حد خوبی نشان دهنده احساسات متن باشد. همچنین صفاتی که معمولا در این چنین متن ها ذکر می شود نیز می تواند راهنمای خوبی برای تشخیص مثبت یا منفی بودن متن باشد (برای مثال کلماتی مانند عالی، افتضاح، زیبا و ... در زبان فارسی).

## تحلیل نتایج رده‌بندی:

با توجه به نتایج استخراج شده می‌توان دید که دقت این روش چیزی در حدود 80 درصد برای تمامی حالات k-fold است که البته نداشتن تفاوت چندان در بین نتایج رده‌بندها قابل پیش‌بینی بود، زیرا در این تمرین ترتیب reviewها تفاوتی در تصمیم‌گیری برای رده‌بند ما ایجاد نمی‌کرد.

pos_precision: 0.815028901734	pos_precision: 0.787234042553	pos_precision: 0.798969072164	pos_precision: 0.788359788359	pos_precision: 0.846590909090
neg_precision: 0.757709251101	neg_precision: 0.726415094339	neg_precision: 0.757281553398	neg_precision: 0.781990521327	neg_precision: 0.78125
pos_recall: 0.719387755102	pos_recall: 0.718446601941	pos_recall: 0.756097560975	pos_recall: 0.764102564102	pos_recall: 0.752525252525
neg_recall: 0.843137254901	neg_recall: 0.793814432989	neg_recall: 0.8	neg_recall: 0.804878048780	neg_recall: 0.866336633663
pos_f_measure: 0.764227642276	pos_f_measure: 0.751269035532	pos_f_measure: 0.776942355889	pos_f_measure: 0.776041666666	pos_f_measure: 0.796791443850
neg_f_measure: 0.798143851508	neg_f_measure: 0.758620689655	neg_f_measure: 0.778054862842	neg_f_measure: 0.793269230769	neg_f_measure: 0.821596244131
accuracy: 0.7825	accuracy: 0.755	accuracy: 0.7775	accuracy: 0.785	accuracy: 0.81

همچنین در مورد موثرترین ویژگی‌های هر رده‌بند می‌توان به نتایج زیر اشاره کرد

ridiculous neg : pos 5.1 : 1.0	waste neg : pos 5.3 : 1.0	awful neg : pos 5.8 : 1.0	waste neg : pos 4.9 : 1.0	awful neg : pos 7.6 : 1.0
awful neg : pos 4.7 : 1.0	ridiculous neg : pos 4.8 : 1.0	waste neg : pos 4.9 : 1.0	awful neg : pos 4.8 : 1.0	ridiculous neg : pos 7.3 : 1.0
waste neg : pos 4.6 : 1.0	awful neg : pos 4.5 : 1.0	ridiculous neg : pos 4.3 : 1.0	ridiculous neg : pos 4.5 : 1.0	waste neg : pos 5.5 : 1.0
worst neg : pos 4.6 : 1.0	worst neg : pos 4.2 : 1.0	dull neg : pos 4.1 : 1.0	worst neg : pos 4.4 : 1.0	worst neg : pos 4.8 : 1.0
memorable pos : neg 4.3 : 1.0	mess neg : pos 4.1 : 1.0	worst neg : pos 3.9 : 1.0	stupid neg : pos 4.3 : 1.0	dull neg : pos 4.6 : 1.0

stupid neg : pos 3.9 : 1.0	dull neg : pos 3.9 : 1.0	subtle pos : neg 3.8 : 1.0	mess neg : pos 3.8 : 1.0	mess neg : pos 4.2 : 1.0
dull neg : pos 3.7 : 1.0	boring neg : pos 3.8 : 1.0	terrible neg : pos 3.7 : 1.0	boring neg : pos 3.8 : 1.0	stupid neg : pos 4.0 : 1.0
subtle pos : neg 3.6 : 1.0	stupid neg : pos 3.8 : 1.0	stupid neg : pos 3.7 : 1.0	memorable pos : neg 3.6 : 1.0	memorable pos : neg 3.6 : 1.0
mess neg : pos 3.4 : 1.0	spielberg pos : neg 3.7 : 1.0	memorable pos : neg 3.6 : 1.0	dull neg : pos 3.5 : 1.0	perfectly pos : neg 3.5 : 1.0
boring neg : pos 3.3 : 1.0	excellent pos : neg 3.7 : 1.0	excellent pos : neg 3.6 : 1.0	realistic pos : neg 3.4 : 1.0	subtle pos : neg 3.4 : 1.0

که به وضوح میزان تاثیرگذاری هر کلمه روی منفی یا مثبت بودن آن نظر مشهود است (علاوه بر کلماتی که به طور کلی صفت‌های مثبت یا منفی هستند، در این لیست کلماتی مانند اسپیلبرگ را مشاهده می‌کنیم که باعث مثبت بودن نظر کاربر نشده و دلیل آن محبوبیت این کارگردان در بین مردم است.)