



دانشکده مهندسی کامپیوتر

مبانی پردازش زبان طبیعی

گزارش پروژه

فاز ۲ - تحلیل، استخراج ویژگی و رده‌بندی

مهدی امیری

۱	لینک به ریپوی پروژه در GitHub	۲
۲	آماده‌سازی محیط مجازی برای اجرای پروژه	۲
۳	آموزش و تحلیل بردار کلمات (word2vec)	۳
۳	۱.۳ شبیه‌ترین کلمات مشترک بین جفت لیبل‌ها	۳
۳	۲.۳ متفاوت‌ترین کلمات مشترک بین جفت لیبل‌ها	۳
۴	tokenization	۴
۴	۱.۴ بخش ۱ داده‌ها	۴
۴	۲.۴ بخش ۲ داده‌ها	۴
۴	۳.۴ بخش ۳ داده‌ها	۴
۴	۴.۴ بخش ۴ داده‌ها	۴
۴	۵.۴ بخش ۵ داده‌ها	۴
۴	۶.۴ انتخاب بهترین تنظیم Tokenizer	۴

۱ لینک به ریپوی پروژه در GitHub

<https://github.com/AmiriShavaki/IUST-NLP-Project>

۲ آماده‌سازی محیط مجازی برای اجرای پروژه

شما می‌توانید به کمک conda و از طریق فایل requirements.txt که در root پروژه در اختیار شما قرار داده شده، به راحتی و فقط با اجرای دستور زیر تمامی پکیج‌های موردنیاز برای اجرای پروژه را نصب کنید: `conda create -name <envname> -file requirements.txt`

۳ آموزش و تحلیل بردار کلمات (word2vec)

در این جا معیار بکارگرفته شده برای شباهت سنجی بین بردارها، شباهت کسینوسی بوده است که از رابطه $\frac{A.B}{|A||B|}$ بدست می آید.

۱.۳ شبیه ترین کلمات مشترک بین جفت لیبل ها

۴ ستاره و ۵ ستاره	۳ ستاره و ۵ ستاره	۳ ستاره و ۴ ستاره	۲ ستاره و ۵ ستاره	۲ ستاره و ۴ ستاره	۲ ستاره و ۳ ستاره	۱ ستاره و ۵ ستاره	۱ ستاره و ۴ ستاره	۱ ستاره و ۳ ستاره	۱ ستاره و ۲ ستاره	کلمه
companions	ither	grew	luxe	twich	refunded	offerings	nighttime	scroller	obstructive	شباهت
0.76	0.90	0.93	0.93	0.97	0.93	0.76	0.79	0.94	0.87	

۲.۳ متفاوت ترین کلمات مشترک بین جفت لیبل ها

۴ ستاره و ۵ ستاره	۳ ستاره و ۵ ستاره	۳ ستاره و ۴ ستاره	۲ ستاره و ۵ ستاره	۲ ستاره و ۴ ستاره	۲ ستاره و ۳ ستاره	۱ ستاره و ۵ ستاره	۱ ستاره و ۴ ستاره	۱ ستاره و ۳ ستاره	۱ ستاره و ۲ ستاره	کلمه
bye	firmware	bye	tnx	reccomended	tnx	science	manages	ha	baffling	شباهت
-0.47	-0.25	-0.23	-0.60	-0.23	-0.34	-0.17	-0.37	-0.41	-0.56	

نکته جالب تفاوت کمتر بردارها در جفت (۱ ستاره، ۵ ستاره) نسبت به بقیه جفت ها است. در این جفت حداکثر تفاوت برداری 0.17 بوده در حالی که همانطور که در جدول می بینید در بقیه جفت ها میزان تفاوت متفاوت ترین بردارها بیشتر است.

کلمه "tnx" در جفت (۵ ستاره و ۲ ستاره) بیشترین تفاوت برداری را داشته است. احتمالاً بخاطر اینکه در داده های ۵ ستاره این کلمه همراه با رضایت مندی و خوشحالی کاربر بیان شده در حالی که در کلاس ۲ ستاره این کلمه توأم با نارضایتی و شاید با لحن کنایه آمیز بیان شده است و تفاوت برداری احتمالاً ناشی از همین تفاوت معنی در بکارگیری کلمه بوده است.

در جفت (۳ ستاره و ۴ ستاره) کلمه ی "reccomended" که دارای غلط املائی است دارای بیشترین تفاوت برداری است. احتمالاً تعداد رخ دادن این کلمه با این غلط املائی خاص در بین داده ها کم بوده و وقتی بطور تصادفی در هر دو کلاس ۳ ستاره و ۴ ستاره از این غلط املائی یافت شده، مدل تعداد زیادی از این کلمه با این شکل نگارشی را در این دو کلاس ندیده که بتواند تشابهات معنایی آن ها را درک کند و احتمالاً در همین تعداد محدودی که رخ داده هم context معنایی متن در دو کلاس متفاوت بوده.

۴ tokenization

۱.۴ بخش ۱ داده‌ها

نتیجه آموزش روی کلاس‌های ۲،۳،۴،۵ ستاره و تست روی کلاس ۱ ستاره:

	vocabsize=50	vocabsize=800	vocabsize=5000	vocabsize=15000
UNK percentage	2.02	4.90	7.29	7.67

۲.۴ بخش ۲ داده‌ها

نتیجه آموزش روی کلاس‌های ۱،۳،۴،۵ ستاره و تست روی کلاس ۲ ستاره:

	vocabsize=50	vocabsize=800	vocabsize=5000	vocabsize=15000
UNK percentage	1.73	4.33	6.44	6.78

۳.۴ بخش ۳ داده‌ها

نتیجه آموزش روی کلاس‌های ۱،۲،۴،۵ ستاره و تست روی کلاس ۳ ستاره:

	vocabsize=50	vocabsize=800	vocabsize=5000	vocabsize=15000
UNK percentage	1.70	4.28	6.32	6.65

۴.۴ بخش ۴ داده‌ها

نتیجه آموزش روی کلاس‌های ۱،۲،۳،۵ ستاره و تست روی کلاس ۴ ستاره:

	vocabsize=50	vocabsize=800	vocabsize=5000	vocabsize=15000
UNK percentage	1.77	4.38	6.56	6.91

۵.۴ بخش ۵ داده‌ها

نتیجه آموزش روی کلاس‌های ۱،۲،۳،۴ ستاره و تست روی کلاس ۵ ستاره:

	vocabsize=50	vocabsize=800	vocabsize=5000	vocabsize=15000
UNK percentage	2.07	4.91	7.41	7.87

۶.۴ انتخاب بهترین تنظیم Tokenizer

طبق خروجی‌هایی که در جدول‌های بالا می‌بینید حالت‌های با vocabsize کوچک‌تر تعداد توکن <UNK> کمتری در هنگام تست ایجاد می‌کرده‌اند. احتمالاً دلیل بهتر عمل کردنشان این است

که هر چقدر vocabsize کوچکتر باشد مدل زیر کلمه‌های کوچکتر و عمومی‌تری را انتخاب می‌کند و در نتیجه در مواجهه با کلمات بیشتر و کلی‌تری از آن‌ها می‌تواند استفاده کند. بهترین عملکرد طبق جداول بالا را مدل آموزش دیده روی کلاس‌های ۱، ۲، ۴، ۵ با vocabsize=50 داشته است پس در فایل run-phase2.bat دستور منتقل کردن این مدل به پوشه models به صورت hard-code اضافه شده است.