

به نام خدا



تمرین سوم درس پردازش زبان طبیعی آماری

«ابهام‌زدایی معنایی کلمات»

استاد درس: دکتر ممتازی

دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر

خرداد ۱۴۰۰



برای ارسال تمرین به نکات زیر توجه کنید.

۱- در جدول زیر نحوه اعمال نمره منفی برای تاخیر در ارسال تمرین‌ها ذکر شده‌است.

میزان جریمه	میزان تاخیر (روز)
هر روز ۵٪	۱ الی ۲ روز
هر روز ۱۰٪	۲ الی ۶ روز

در صورتی که برای ارسال تمرین‌ها بین ۷ تا ۱۴ روز تاخیر داشته باشید، نمره شما از ۵۰٪ محاسبه می‌شود و پس

از این بازه به تمرین ارسالی نمره‌ای تعلق نمی‌گیرد.

۲- هرگونه کپی‌برداری در انجام تمرین‌ها موجب کسر نمره خواهد شد.

۳- آخرین مهلت ارسال تمرین، ساعت ۲۳:۵۵ روز دوشنبه ۳۱ خرداد می‌باشد.

۴- فایل‌های ارسالی خود شامل فایل‌های پیاده‌سازی و گزارش را فشرده کنید و با عنوان «شماره دانشجویی_HW3»

مانند HW3_97131022 ارسال کنید.

۵- زبان برنامه‌نویسی برای انجام تمرین‌ها، پایتون یا جاوا در نظر گرفته شده‌است.

۶- کدهای ارسالی خود را برای افزایش خوانایی و درک بهتر به صورت مناسب کامنت‌گذاری کنید.

۷- در صورت هرگونه سوال یا مشکل می‌توانید با تدریس‌یار درس از طریق ایمیل زیر در ارتباط باشید.

mousavian12@gmail.com

مریم موسویان

بخش اول: تعریف مسئله و معرفی دادگان

در این تمرین قصد داریم وظیفه ابهام‌زدایی معنایی کلمات^۱ را به عنوان یک مولفه اساسی و مهم در کاربردهای مختلف پردازش زبان طبیعی با استفاده از الگوریتم‌های یادگیری ماشین انجام دهیم. مجموعه دادگان و کلیه فایل‌های مورد نیاز برای این تمرین از لینک^۲ موجود در پاورقی قابل دانلود است. در فایل train.csv ۶۰۹۸ داده آموزش و در فایل test.csv ۱۱۰۶ داده آزمون قرار دارد. مجموعه کلمات و تعداد نمونه‌های آن‌ها در هر فایل در جدول زیر مشاهده می‌شوند. ویژگی‌های موجود در مجموعه داده به ترتیب کلمه (word)، اجزای سخن (pos^۳)، شناسه نمونه (instance_id)، منبع سند (doc_src)، شناسه مفهوم (sense_id) و بافت کلمه (context) است. (کلمه مبهم با تگ <head> درون context مشخص شده‌است). معنی هر مفهوم را با توجه به شناسه آن مفهوم و کلمه مورد نظر در فایل dictionary.xml می‌توانید بیابید.

داده آزمون	داده آموزش	کلمه	داده آزمون	داده آموزش	کلمه	داده آزمون	داده آموزش	کلمه
17	92	organization	27	147	disc	28	156	activate
28	152	paper	26	146	eat	35	192	add
32	181	party	20	110	encounter	38	214	appear
24	132	performance	21	114	expect	26	146	argument
23	124	plan	17	92	express	40	223	arm
16	88	play	10	53	hear	31	175	ask
28	156	produce	13	71	hot	19	104	atmosphere
21	115	provide	22	122	image	26	143	audience
8	44	receive	4	22	important	38	213	bank
21	118	remain	28	154	interest	24	133	begin
9	50	rule	10	52	judgment	20	110	climb
25	141	shelter	11	60	lose	19	103	decide
5	27	simple	13	67	mean	34	191	degree
16	88	smell	9	49	miss	30	169	difference
8	45	solid	20	110	note	14	78	different
23	127	sort	6	29	operate	7	35	difficulty
16	84	watch	17	94	treat	9	46	source
12	65	win	4	22	use	19	107	suspend
10	37	write	10	56	wash	22	124	talk

¹ Word Sense Disambiguation (WSD)

² <https://drive.google.com/drive/folders/1FidQTKvT4YAyKD4ncqIizjVE4POs7VAB?usp=sharing>

³ Part Of Speech

توجه:

- نحوه انجام پیش‌پردازش بر روی داده‌ها شامل کتابخانه مورد استفاده و مراحل انجام‌شده را در گزارش خود مکتوب کنید.
- برای انجام این تمرین می‌توانید از کتابخانه‌های آماده استفاده کنید.

بخش دوم: انتخاب ویژگی

در این قسمت با استفاده از مجموعه داده‌گان ذکرشده در بخش قبل، بازنمایی برای کلمات دارای ابهام معنایی در مجموعه داده به دست آورید. در این تمرین می‌خواهیم از مدل بازنمایی مبتنی بر بافت⁴ BERT و مدل Word2Vec برای استخراج ویژگی کمک بگیریم. معمولاً برای استخراج بازنمایی کلمات از مدل BERT، وزن‌ها و پارامترهای مدل را روی پیکره مورد نظر به صورت دقیق تنظیم می‌کنند. اما در این درس به دلیل عدم آشنایی برخی دانشجویان با درس شبکه‌های عصبی، از تنظیم وزن‌ها و پارامترهای مدل صرف نظر می‌شود. بنابراین برای استخراج بازنمایی‌های مبتنی بر بافت BERT می‌توانید از کتابخانه bert-embedding⁵ استفاده کنید (استفاده از سایر کتابخانه‌ها نیز مانعی ندارد). خروجی تمام بردارها ۳۰۰ بعد در نظر گرفته شود. (با توجه به اینکه بردار خروجی مدل BERT برای هر کلمه ۷۶۸ بعد است و برای این که در هنگام آموزش دسته‌بندها در قسمت بعد با کمبود منابع رو به رو نشوید، بردارهای خروجی برای کلمات مبهم را با استفاده از الگوریتم PCA به ۳۰۰ بعد کاهش دهید.) برای استخراج بازنمایی کلمات از مدل Word2Vec از مدل از پیش آموزش دیده گوگل بر روی مجموعه داده Google News استفاده کنید. (این مدل در لینک حاوی داده‌ها و فایل‌های تمرین با عنوان GoogleNews-vectors-negative300.bin.gz آپلود شده است. برای دانلود و لود کردن این مدل نیز می‌توانید از کتابخانه gensim استفاده کنید.)

با استفاده از روش‌های زیر بازنمایی هریک از ورودی‌های مبهم را به دست آورید و سپس با استفاده از بازنمایی و توضیحات بخش بعد مفهوم مرتبط با آن ورودی را به دست آورید.

الف) استفاده از مدل از پیش آموزش داده شده BERT و استخراج بازنمایی تنها برای کلمه مبهم موجود در متن از مدل BERT بدون استخراج بازنمایی کلمات دیگر متن. (با توجه به اینکه مدل BERT مبتنی بر بافت است کلمات اطراف کلمه هدف تاثیر خود را در بازنمایی کلمه هدف به جا می‌گذارند اما نیاز به بازنمایی آن‌ها به صورت مجزا نیست.)

⁴ Bidirectional Encoder Representations from Transformers (BERT)

⁵ <https://pypi.org/project/bert-embedding>

(ب) استفاده از مدل از پیش آموزش داده شده BERT و استخراج بازنمایی کلمات از مدل BERT برای کلمه هدف و سایر کلمات موجود در بافت کلمه هدف به شعاع ± 3 کلمه و سپس استفاده از میانگین وزن دار بازنمایی کلمات موجود در بافت کلمه مبهم با استفاده از TF-IDF هر یک از کلمات. (منظور از میانگین وزن دار این است که TF-IDF هر کلمه به عنوان وزن بردار بازنمایی آن کلمه در نظر گرفته شود).

(ج) تکرار بند (ب) ولی با در نظر گرفتن تمام کلمات موجود در بافت کلمه مبهم.

(د) استفاده از بازنمایی از پیش آموزش داده شده Word2Vec برای کلمه هدف و سایر کلمات موجود در بافت کلمه هدف به شعاع ± 3 کلمه و سپس استفاده از میانگین وزن دار بازنمایی کلمات موجود در بافت کلمه مبهم با استفاده از TF-IDF هر یک از کلمات.

(ه) تکرار بند (د) ولی با در نظر گرفتن تمام کلمات موجود در بافت کلمه مبهم.

بخش سوم: دسته بندی

در این قسمت می خواهیم بهترین مفهوم از مجموعه مفاهیم داده شده برای کلمه مبهم را با توجه به بافت به کار رفته برای آن کلمه بیابیم. برای این کار از الگوریتم های یادگیری ماشین^۶ و روش های یادگیری گروهی^۷ استفاده می کنیم. (انتخاب بهترین مقدار برای پارامترهای دسته بندی های استفاده شده به عهده شما می باشد).

(الف) برای دسته بندی از روش های بازنمایی مطرح شده در بخش قبل و دسته بندی های Logistic Regression و Random Forest استفاده کنید. نتایج آزمایش های ذکر شده را با معیارهای Accuracy و F-Measure ارزیابی کنید. معیارهای ارزیابی را برای اسم ها، فعل ها و صفت ها به صورت مجزا، میانگین نتایج را گزارش کنید.

(ب) با توجه به ۵ بازنمایی مختلف و دو دسته بند مختلف، ۱۰ حالت را در این تمرین آزمایش نموده اید. از میان این ۱۰ حالت ۳ نتیجه برتر را انتخاب کنید و با کمک این سه مدل یک مدل گروهی^۸ بسازید. برای پیدا کردن مفهوم مناسب برای هر کلمه مبهم بین نتایج ۳ دسته بند رای اکثریت بگیرید. مانند قسمت قبل معیارهای ارزیابی Accuracy و F-Measure را برای هر POS به صورت جداگانه گزارش کنید.

(ج) با استفاده از ۳ تا بهترین دسته بند انتخاب شده در بند (ب) یک مدل گروهی دیگر بسازید. برای این کار یک دسته بند دلخواه را با استفاده از خروجی ۳ دسته بند (مفاهیم پیش بینی شده برای کلمات مبهم) آموزش دهید. در واقع ویژگی های استفاده شده

⁶ Machine Learning

⁷ Ensemble Learning

⁸ Ensemble Model

برای آموزش دسته‌بند دلخواه یک بردار با ابعاد ۳ شامل برچسب‌های (مفاهیم) پیش‌بینی شده توسط ۳ دسته‌بند خواهد بود. مانند قسمت‌های قبل معیارهای ارزیابی Accuracy و F-Measure را برای هر POS به صورت جداگانه گزارش کنید. (د) در این قسمت می‌خواهیم برای ابهام‌زدایی معنایی کلمات از شبکه‌های عصبی استفاده کنیم. ابتدا بازنمایی توکن [cls] را از مدل BERT استخراج کنید و به یک شبکه عصبی پیش‌خور^۹ دهید (برای استفاده از مدل BERT به ابتدای هر متن ورودی، توکن [cls] اضافه می‌شود. بازنمایی که از مدل BERT برای این توکن در هر متن استخراج می‌شود، به عنوان بازنمایی برای کل متن در نظر و از آن در وظیفه دسته‌بندی متون کمک گرفته می‌شود). مانند قسمت‌های قبل معیارهای ارزیابی Accuracy و F-Measure را برای هر POS به صورت جداگانه گزارش کنید. (این قسمت اختیاری و دارای نمره اضافه خواهد بود).

بخش چهارم: تحلیل نتایج

در این قسمت نتایج به دست آمده از دسته‌بندهای مختلف بخش سوم را مقایسه و به صورت کامل تحلیل کنید. در صورتی که با استفاده از ایده‌ای خلاقانه بتوانید نتایج را بهبود دهید نمره‌ای اضافه دریافت خواهید کرد.

موفق باشید

⁹ Feed Forward Neural Network