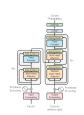


تمرین صفرم پیشپردازش زبان طبیعی

تاریخ بارگزاری تمرین: 13/12/1401 تاریخ تحویل تمرین: 22/12/1401



۱ پرسشها

فرآیند ریشهیابی ابا لمسازی کم تفاوتهایی دارد؟ مثال بزنید.

۲ برنامەنوپسى

دستهبندی یکی از عملیاتهای مرسوم یادگیری ماشین به حساب می آید. یکی از جنبههای این عملیات در زمینه پردازش زبان طبیعی میشندی متون است. در فایلهای مربوط به تمرین، کد و مجموعه دادگان که برنامه ید دستهبندی ساده آمده است که براساس الگوریتم دستهبندی kنزدیک ترین همسایگی می می کند. دادگان موجود مجموعه ای از نظرات سایت یوتیوب می تشکیل می دهند. هدف این دستهبند، تشخیص اسپم بودن یا نبودن این نظرات است. بدون اعمال پیش پردازش روی دادگان، این دسته بند به صحت 8 درصد می رسد. مجموعه ای از عملیاتهای پیش پردازش را به گونه ای انتخاب کنید که با اعمال آنها روی دادگان، صحت مدل یاد شده افزایش یابد. در این راستا موارد زیر را در نظر بگیرید:

- ۱. گام نهایی در عملیات پیشپردازش میبایست یکی از الگوریتمهای ریشهیابی، لمسازی و تصحیح غلط۳۱ باشد.
- ۲. عملیات دستهبندی بایستی بهازای هر کدام از الگوریتمهای یاد شده، یکبار و با مجموعه اَعمال پیشیردازش یکسان اجرا گردد.
- ۳. عملیاتهای پیشپردازش بایستی باتوجه به هدف مسئلهی دستهبندی و کثیفیهای مجموعهدادگان طراحی گردند.
 - ۴. حداقل تعداد پیشیردازش ممکن بایستی استفاده گردد.
 - ۵. حداقل بهبود نتیجه صحت 8درصد خواهد بود.

¹Stemming

²Lemmatization

³Classification

⁴Machine learning

⁵Natural language processing

⁶Code

⁷Dataset

⁸K nearest neighbours

⁹Comments

¹⁰Youtube

¹¹Spam

¹²Accuracy

¹³Spell correction

پس از انجام اجراهای خواسته شده مشخص کنید کدام یک از الگوریتمهای ریشهیابی، لمسازی و تصحیح غلط بهتر عمل میکنند. آیا با شخصیسازی روند پیشپردازش برای هر یک از این الگوریتمها میتوان به نتایج بهتری نتیجهی بهتری رسید؟ توضیح دهید. آیا میتوان با جایگزین کردن الگوریتم دستهبندی به نتایج بهتری دست پیدا کرد، الگوریتم پیشنهادی خود را مشخص نمایید و نتایج جدید را گزارش کنید.

٣ خوانش مقاله

مقالهی «On the effectiveness of preprocessing methods when dealing with different levels» مقالهی «of class imbalance را مورد مطالعه قرار دهید و به پرسشهای زیر پاسخ دهید:

- ۱. هدف مقاله را در یک یا دو جمله توضیح دهید.
- ۲. تفاوت این مقاله با مقالات پیش از خود در چیست؟
- ۳. روشهای بیشنمونهبرداری^{۱۱} و کمنمونهبرداری^{۱۵} که در مقاله آمده است را توضیح دهید. از هر کدام یک مثال بیاورید و به صورت خلاصه بگویید چگونه عمل میکنند.
 - ۴. چرا معیار ۱۶ صحت در مورد مجموعهدادگان با عدمتوازن ۱۱^{۱۷}، معیار خوبی نیست؟
- ۵. چه معیارهایی برای بررسی نتایج الگوریتمها روی مجموعهدادگان با عدمتوازن معرفی شده است،
 آنها را به طور خلاصه توضیح دهید؟ چرا بهتر هستند؟
 - آیا مجموعهداده فراهم شده در بخش ۲ دارای عدمتوازن است؟
 - ۷. یکی از معیارهای معرفی شده را به برنامه بخش ۲ اضافه کنید.
 - ٨. نتيجهي مقاله چيست؟

۴ نکات تحویل

- ۱. پاسخ خود را تحت پوشهای به اسم NLP_NAME_ID و در قالب zip بارگذاری نمایید.
 - ۲. این پوشه میبایست حاوی موارد زیر باشد:
 - پوشهای با نام code باشد که شامل برنامهی نوشته/تغییر داده شده است.
 - پوشهای با نام doc که حاوی داکیومنتها و فایل توضیحات میباشد.
 - ۳. لازم به ذکر است که رعایت قوانین نگارشی حائز اهمیت خواهد بود.

مراجع

[1] V. García, J. S. Sánchez, and R. A. Mollineda, "On the effectiveness of preprocessing methods when dealing with different levels of class imbalance," *Knowledge-Based Systems*, vol.25, no.1, pp.13–21, 2012.

¹⁴Over-sampling

¹⁵Under-sampling

¹⁶Metric

¹⁷Imbalance