

به نام آنکه آموخت انسان را آنچه نمردانست



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



درس پردازش زبان‌های طبیعی

CA1

m.ghoroobi@ut.ac.ir

محمد غروبى

بهمن ۱۳۹۸

فهرست

۳	مقدمه
۴	دادگان
۵	سوال ۱ - پیش پردازش (۱۰ نمره)
۶	سوال ۲ - ایجاد مدل های زبانی (۳۰ نمره)
۷	سوال ۳ - معیار سرگشتگی (۴۵ نمره)
۷	الف)
۷	ب)
۹	سوال ۴ - دادگان تست (۱۵ نمره)

هدف از انجام این تمرین، آشنایی با مباحث ان گرام‌ها و معیار سرگشتگی^۲ در زبان فارسی است. انتظار می‌رود در طی این تمرین، با توجه به ان گرام‌ها، مدل‌های زبانی^۳ مربوطه در زبان فارسی ایجاد نمایید، پیش‌بینی‌های خود را با استفاده از معیار سرگشتگی انجام دهید و در نهایت ارزیابی نتایج را با استفاده از معیار F۱ گزارش کنید.

نکته: هدف از تمرین، ایجاد تفکر علمی در شما می‌باشد، پس فارغ از نتایج، بخشی اعظمی از نمره به تفسیر نتایج تعلق می‌گیرد.

^۱ N-grams

^۲ Perplexity

^۳ Language Models

دادگان فارسی از اخبار روزنامه همشهری بین سال‌های ۱۳۷۵ تا ۱۳۸۶ با استفاده از خزشگرهای وب ایجاد

شده است. با توجه به حجم بالای دادگان اصلی، نمونه‌ای از آن ایجاد و در اختیار شما قرار گرفته است.

تعداد ۲۳۸۱ متن خبری در ۶ کلاس (تکنولوژی، ورزشی، اجتماعی، سیاسی، اقتصادی و فرهنگی) در بخش دادگان آموزشی^۴ قرار گرفته است. هم چنین ۶۰۰ متن خبری بدون کلاس در بخش دادگان تست^۵ قرار گرفته است.

Train data^۴

Test data^۵

سوال ۱ – پیش پردازش (۱۰ نمره)

با توجه به کتابخانه‌های آماده در زبان فارسی و کتابخانه‌های آماده دیگر که در اختیار شما قرار گرفته، می‌توانید عملیاتی همچون نشانه‌گذاری، ریشه‌یابی و ... را انجام دهید.

دقت کنید که پیش‌پردازش یکی از مهمترین مراحل پروژه‌های پردازش زبان طبیعی است که کیفیت آن بر روی نتیجه، تاثیر مستقیم دارد. همین‌طور به پیش‌پردازشهایی که خاص زبان فارسی است دقت کنید.

راهنمایی:

منظور از نشانه‌گذاری مشخص نمودن ابتدا و انتهای تمامی جملات است.

علی به مدرسه رفت. بابک آمد.

<s>علی به مدرسه رفت**</s>** **<s>**بابک آمد**</s>**

سوال ۲- ایجاد مدل‌های زبانی (۳۰ نمره)

ابتدا دادگان آموزشی را به دو بخشی آموزشی (۸۰٪) و اعتبارسنجی^۶ (۲۰٪) تقسیم کنید.

با استفاده از دادگان آموزشی و با استفاده از [کتابخانه‌های موجود](#) برای هر کلاس مدل‌های زبانی یکتایی^۷ و دوتایی^۸ برای آن گرام‌های کلمه و حرف ایجاد کنید.

راهنمایی:

مدل یکتایی کلمه: هر کلمه در متن را جداگانه در نظر بگیرید.

مدل دوتایی کلمه: دو کلمه پشت سر هم در متن را در نظر بگیرید.

مدل یکتایی حرف: هر حرف در متن را در نظر بگیرید. (به فاصله دقت شود)

مدل دوتایی حرف: هر دو حرف پشت سر هم در متن را در نظر بگیرید. (به فاصله دقت شود)

۴ مدل که هر یک برای ۶ کلاس ساخته شده است، پس نهایتاً ۲۴ مدل زبانی متفاوت ایجاد می‌گردد.

^۶ Validation

^۷ Unigram

^۸ Bigram

سوال ۳ – معیار سرگشتگی (۴۵ نمره)

الف)

در سوال ۲ مدل‌های زبانی را ایجاد کردید. حال برای هریک از متن‌های خبری در دادگان اعتبارسنجی کلاس مورد نظر را با توجه به حداقل میزان سرگشتگی که از مدل زبانی مربوط بدست می‌آید تعیین کنید.

راهنمایی:

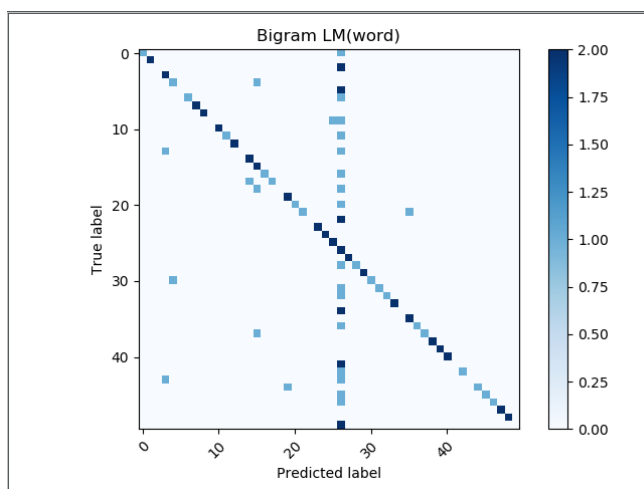
در مدل یکتایی کلمه، برای هر متن خبری، ۶ مدل زبانی دارید، پس ۶ مقدار سرگشتگی بدست می‌آورید، کمترین مقدار سرگشتگی متعلق به هر کلاسی باشد، کلاس متن خبری همان کلاس می‌باشد. حال برای هر متن خبری در داده‌های اعتبارسنجی یک کلاس حقیقی (از پیش تعیین شده) دارید و یک کلاس تخمینی که با توجه به معیار سرگشتگی بدست آورده‌اید، (پس براحتی به معیارهای دقت، صحت و درنتیجه ۱f دسترسی دارید) برای بدست آوردن سرگشتگی بر روی مدل‌های زبانی می‌توانید از [کتابخانه‌های موجود](#) کمک بگیرید.

ب)

در نهایت ۴ مدل زبانی داریم، که برای هر یک، یک ماتریس درهم‌ریختگی^۹ ایجاد کنید و برای هر نمودار تنها معیار میانگین ۱f تمامی کلاس‌ها را گزارش کنید.

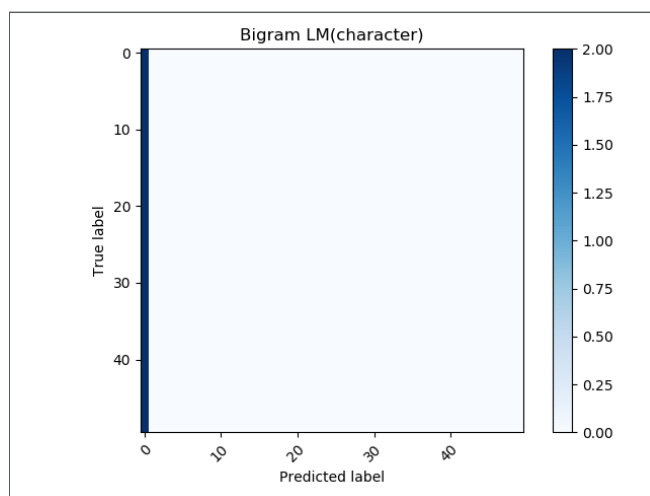
راهنمایی:

نمونه‌ای از ماتریس درهم‌ریختگی با معیار ۱f بالا



Confusion Matrix^۹

نمونه‌ای از ماتریس درهم ریختگی با معیار f_1 پایین



سوال ۴- دادگان تست (۱۵ نمره)

جهت ارزیابی نتایج شما، دادگان تست که فاقد کلاس می باشند در اختیار شما قرار گرفته. با توجه به بهترین مدل زبانی که در سوال ۳ بدست آورید. لطفا برچسب هر متن به فرمت زیر در فایلی تحت عنوان Result.csv ارسال کنید.

Filename ,Class

۱۰۰۳۸.txt ,finance

به قواعد حل تمرین در سایت حتما توجه گردد.

موفق باشید.