

به نام خدا



تمرین اول درس پردازش زبان طبیعی
«آشنایی با مدل‌های زبانی و روش‌های بازنمایی برداری کلمات»

استاد درس: دکتر ممتازی

دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر

زمستان ۱۴۰۲



برای ارسال تمرین به نکات زیر توجه کنید.

- ۱- برای ارسال پاسخ تمرین‌های این درس، **مجموعاً ۱۰ روز** زمان تاخیر مجاز در نظر گرفته شده است و در صورتی که مجموع زمان تاخیرها از این مقدار بیشتر باشد، پاسخ ارسال شده مورد بررسی قرار نخواهد گرفت.
- ۲- هرگونه کپی کردن در انجام تمرین‌ها موجب کسر نمره خواهد شد.
- ۳- آخرین مهلت ارسال تمرین، **ساعت ۲۳:۵۵ روز شنبه ۱۸ فروردین** می‌باشد.
- ۴- فایل‌های ارسالی جهت نمره‌دهی باید شامل پیاده‌سازی و گزارش تمرین باشد، فایل‌های خود را فشرده نمایید و به صورت «**شماره دانشجویی_HW1**» مانند **HW1_400131022** نام‌گذاری کنید.
- ۵- زبان برنامه‌نویسی برای انجام این تمرین، تنها می‌تواند **پایتون** باشد.
- ۶- کدهای ارسالی خود را برای افزایش خوانایی و درک بهتر به صورت مناسب کامنت‌گذاری کنید.
- ۷- در این تمرین شما باید **موارد خواسته شده را پیاده‌سازی نمایید** و استفاده از کتابخانه‌های آماده مجاز نمی‌باشد.
- ۸- در صورت هرگونه سوال یا مشکل می‌توانید با تدریس‌یار درس از طریق ایمیل زیر در ارتباط باشید.

mohammad.naeimi+nlp@aut.ac.ir

محمد نعیمی

بخش اول: تعریف مسئله و معرفی مجموعه داده

مجموعه داده در این تمرین، شامل ۳ فایل `train.csv`، `val.csv` و `test.csv` که به ترتیب برای آموزش^۱، ارزیابی^۲ و آزمون^۳ می باشد، در اختیار شما قرار گرفته است. این مجموعه داده مربوط به دسته بندی متن می باشد که مجموعه ای از مقالات خبری فارسی است. هدف ما در این تمرین این است که با استفاده از این مجموعه داده مدل زبانی احتمالاتی تولید کنیم. همچنین بازنمایی برداری هر کلمه را تولید نموده، برای تعدادی از کلمات، کلمات هم معنی را پیدا کنیم. مجموعه داده ارائه شده شامل دو ستون زیر می باشد:

ویژگی	توضیحات
content	متن داده
label	برچسب داده

همچنین تعداد محتوای هریک از فایل ها به شرح زیر می باشد:

نام فایل	تعداد ورودی ها
train	13,314
val	1,480
test	1,644

بخش دوم: مدل های زبانی احتمالاتی (۴۰ امتیاز)

در این قسمت سه مدل زبانی احتمالاتی یونیگرم^۴، بایگرم^۵ و ترايگرم^۶ را بر روی مجموعه داده آموزش ایجاد کنید. هریک از مدل های زبانی ایجاد شده را به صورت مجزا به دو روش `Back-off Smoothing` و `Absolute Discounting` هموارسازی^۷ نموده، مقدار بهینه پارامترهای هموارسازی را با استفاده از مجموعه داده ارزیابی محاسبه کنید. برای بررسی عملکرد مدل های زبانی ایجاد شده معیار ارزیابی `perplexity` را بر روی مجموعه داده آزمون محاسبه کنید.

الف) معیار ارزیابی را برای تمامی متن های مجموعه داده آزمون گزارش کنید.

ب) معیار ارزیابی را به صورت مجزا برای هر برچسب مجموعه داده آزمون گزارش کنید.

ج) نتایج مدل های مختلف را با یکدیگر مقایسه نموده، تاثیر هموارسازی های متفاوت را بررسی و تحلیل کنید.

¹ Train
² Validation
³ Test
⁴ Unigram
⁵ Bigram
⁶ Trigram
⁷ Smoothing

بخش سوم: روش‌های بازنمایی برداری کلمات (۴۰ امتیاز)

می‌خواهیم با استفاده از مجموعه داده ارائه شده، بازنمایی برداری متن‌ها را تولید نموده و آن‌ها را دسته‌بندی کنیم. بعد از تولید بردار متناظر با هر متن، یک مدل دسته‌بندی را با استفاده از مجموعه داده آموزش، آموزش داده و با استفاده از مجموعه داده آزمون، ارزیابی کنید. می‌توانید از هر روش دلخواه (مانند KNN یا شبکه‌های عصبی) برای دسته‌بندی متن‌ها استفاده کنید. توجه کنید که خروجی بردار متن حداقل ۲۰۰ در نظر گرفته شود. برای تولید بردارهای متن‌ها، هر یک از روش‌های زیر را به صورت جداگانه استفاده کنید، نتایج معیارهای ارزیابی دسته‌بندی را گزارش نموده، مقایسه و تحلیل نمایید. (برای استفاده از مدل word2vec، کتابخانه genism پیشنهاد می‌شود).

(ب) میانگین حسابی بردارهای word2vec (در حالت skip-gram) کلمات متن.

(ج) میانگین وزنی بردارهای word2vec (در حالت skip-gram) کلمات متن، با مقدار TF-IDF هر کلمه به عنوان وزن.

بخش چهارم: استفاده از روش‌های ارزیابی و تحلیل نتایج (۲۰ امتیاز)

الف) همانطور که توضیح داده شد، برای بررسی عملکرد مدل‌های زبانی ایجاد شده، معیار ارزیابی perplexity را بر روی مجموعه داده آزمون محاسبه می‌کنیم. این معیار ارزیابی را پیاده‌سازی نمایید. توضیح دهید perplexity چیست و چگونه برای ارزیابی مدل‌های زبانی استفاده و تفسیر می‌شود؟

(ب) معیارهای ارزیابی Accuracy و F1-score را برای دسته‌بندی چندکلاسه پیاده‌سازی و استفاده نمایید.

موفق باشید

محمد نعیمی