

	<p>تمرین اول – درس پردازش زبان طبیعی آماری دکتر ممتازی زمستان ۹۷ – دانشکده مهندسی کامپیوتر، دانشگاه صنعتی امیرکبیر</p>
---	--

1. در جدول زیر نحوه اعمال نمره منفی برای تاخیر در ارسال تمرین ها آورده شده است:

میزان نمره منفی	تاخیر (روز)
هر روز 5٪	از 1 الی 2
هر روز 10٪	از 3 الی 6

توجه داشته باشید در صورت تاخیر بین 7 تا 14 روز، نمره تمرین از 50٪ محاسبه شده و پس از آن نمره ای تعلق نمی گیرد.

2. هدف از انجام تمرین ها یادگیری عمیق تر مطالب درسی است. در نتیجه هرگونه کپی برداری موجب کسر نمره خواهد شد.

3. تا ساعت 23:15 روز جمعه ۱۶ فروردین فرصت دارید تمرین را در مودل بارگذاری کنید. تمام فایل های پیاده سازی را به همراه فایل pdf مربوط به گزارش تمرین را در یک فایل فشرده قرار دهید. نام فایل نهایی را شماره دانشجویی خود (برای مثال HW1_97131024) قرار دهید.

4. زبان برنامه نویسی برای انجام تمرین: پایتون، جاوا و یا متلب در نظر گرفته شده است.

5. برنامه های نوشته شده خوانا باشد و کامنت گذاری مناسب باشد (طوری که روند کار کاملاً مشخص باشد).

6. در صورت وجود هر گونه سوال می توانید از طریق ایمیل با تدریس یاران درس در ارتباط باشید:

javadforough@gmail.com و amin.ghsm@aut.ac.ir

مساله: دسته بندی^۱متون

بخش اول: استخراج ویژگی^۲

در این تمرین قصد داریم عمل دسته بندی متون را به روش بیز با استفاده از مدل های Unigram و Bigram انجام دهیم. برای انجام این کار بایستی از کلیه کلمات موجود در پیکره آموزشی به عنوان ویژگی استفاده شود. اما به دلیل حجم محاسبات مورد نیاز ما می خواهیم یک مرحله انتخاب ویژگی پیش از انجام دسته بندی صورت گیرد.

بدین منظور، پیکره^۳ کوتاه شده همشهری در اختیار شما قرار می گیرد. این پیکره در یک فایل تدوین شده است که هر خط فایل یک سند^۴ را شامل می شود. در ابتدای هر خط نیز برچسب آن سند نوشته شده است و با علامت @ از متن جدا شده است.

در این راستا می خواهیم تنها برای حالت unigram انتخاب ویژگی با استفاده از الگوریتم information gain صورت گیرد. خروجی مورد انتظار به منظور گزارش:

- ۲۰۰ عدد از بهترین ویژگی های بدست آمده به همراه امتیاز مربوط به هر یک از ویژگی ها

بخش دوم: دسته بندی

در این بخش برای حالت unigram دسته بندی را به دو صورت زیر انجام دهید:

الف) با استفاده از ویژگی های بدست آمده در مرحله ی قبل

ب) عدم استفاده از انتخاب ویژگی

همچنین برای حالت bigram بدون انجام انتخاب ویژگی عمل دسته بندی صورت گیرد.

لازم به ذکر است که در این مرحله بایستی از هموارسازی^۵ به روش الگوریتم Absolute Discounting به منظور حل مشکل الگوهای دیده نشده^۶ استفاده شود. این هموارسازی را با استفاده از مقادیر 0.1 ، 0.3 و 0.5 برای δ انجام دهید.

در این بخش شما لازم هست موارد زیر را در گزارش مربوطه ارایه نمایید:

- ماتریس سردرگمی^۷ مربوط به تست انجام شده با دسته بند بیز با استفاده از مدل های Unigram و Bigram (همزمان با اعمال Absolute Discounting برای هر سه مقدار δ)
- بررسی تأثیر مقادیر مختلف δ بر روی عملکرد مدل ها

¹Classification

²Feature

³Bayes

⁴Corpus

⁵Document

⁶Label

⁷Smoothing

⁸Unseen Patterns

⁹Confusion Matrix

- ارزیابی و تحلیل دسته بندی های انجام شده با استفاده از معیارهای زیر برای هر دو دسته بند

- Precision

- Recall

- F1-measure

- Macro Average و Micro Average برای هر سه معیار بالا

لازم به ذکر است که لینک پیکره همشهری (شامل داده آموزش و داده آزمون) در زیر آورده شده است:

<https://www.dropbox.com/s/zkmozli0mmentpy/HAM-Train-Test.zip?dl=0>

موفق و پیروز باشید