

	<p>تمرین دوم – درس پردازش زبان طبیعی آماری</p> <p>دکتر ممتازی</p> <p>بهار ۹۸ – 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 ۰.۵٪	از ۱ الی ۲
هر روز ۱.۰٪	از ۳ الی ۶

توجه داشته باشید در صورت تاخیر بین ۷ تا ۱۴ روز، نمره تمرین از ۵۰٪ 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 ۲۳:۵۵ روز دوشنبه ۲۳ اردیبهشت 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 HW2\_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:

[amin.ghsm@aut.ac.ir](mailto:amin.ghsm@aut.ac.ir) و [javadforough@gmail.com](mailto:javadforough@gmail.com)

## بازنمایی و خوشه بندی متن

در این تمرین می خواهیم بازنمایی متن را به چندین روش محاسبه نموده و سپس روی این بازنمایی ها با استفاده از الگوریتم kmeans به خوشه بندی متن بپردازیم. برای این منظور همانند تمرین اول بخشی از مجموعه دادگان همشهری در نظر گرفته شده است. در این راستا دو بخش را بایستی انجام دهید.

### بخش اول: بازنمایی

در این قسمت با استفاده از پیکره ی ارائه شده، بازنمایی متن ها را بدست می آوریم. برای این کار ۴ روش زیر را لحاظ کنید.

- آموزش بردار کلمات روی پیکره ارائه شده (با استفاده از Word2Vec مدل Skip-gram) و سپس استفاده از میانگین بازنمایی کلمات متن به منظور محاسبه بازنمایی متن
- آموزش بردار کلمات روی پیکره ارائه شده (با استفاده از Word2Vec مدل Skip-gram) و سپس استفاده از میانگین وزن دار بازنمایی کلمات متن (با به کار بستن tf-idf هر یک از کلمات) به منظور محاسبه بازنمایی متن
- آموزش بردار متون با استفاده از مدل doc2vec روی پیکره ی ارائه شده
- ساخت ماتریس doc-word با استفاده از اطلاعات tf کلمات و سپس کاهش بعد ماتریس از طریق روش SVD

### بخش دوم: خوشه بندی

در این بخش با استفاده از بازنمایی های بدست آمده در بخش قبل الگوریتم خوشه بندی Kmeans را اجرا کرده و پارامتر تعداد خوشه را برابر ۶ در نظر بگیرید. بدیهی است که هنگام خوشه بندی هیچ گونه استفاده از از برچسب اسناد صورت نخواهد گرفت. حال برای ارزیابی خوشه های ایجاد شده دو سناریوی ارزیابی زیر را اجرا نمایید. یکی ارزیابی بر روی مجموعه داده آموزش و دیگری ارزیابی بر روی داده آزمون که توضیح آن به شرح زیر است:

#### ❖ ارزیابی آموزش

بر چسب هر خوشه را برابر پرتکرارترین برچسب آن خوشه در نظر بگیرید. در نتیجه تمامی نمونه ها در یک خوشه، پرتکرارترین برچسب را می گیرند که به عنوان برچسب حدس زده شده توسط مدل در نظر گرفته می شود. همچنین برچسب های واقعی نیز همان برچسب های مشخص شده با @ در مجموعه آموزشی می باشد. با داشتن برچسب گلد و برچسب حدس زده شده کیفیت خوشه بندی را با استفاده از ۴ معیار Accuracy, NMI, F-Measure و V-Measure بدست آورید.

#### ❖ ارزیابی آزمون

هر یک از نمونه های آزمون را با مراکز خوشه مقایسه کرده و هر کدام که نزدیک تر بود برچسب آن خوشه (که در ارزیابی آموزش توضیح داده شد) به عنوان برچسب حدس زده شده توسط مدل خوشه بندی در نظر گرفته می شود. مجدداً کیفیت خوشه بندی را با استفاده از ۴ معیار Accuracy, NMI, F-Measure و V-Measure بدست آورده و نتایج را بررسی نمایید.

### توجه:

- برای انجام تمرین می توانید از کتابخانه های آماده نظیر Gensim, NLTK, Scikit و ... استفاده نمایید.
- برای تمام روش های بازنمایی خروجی بردار هر متن ۳۰۰ بعد باشد
- 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>