

درس كلان دادهها

تمرین دوم



نیم سال دوم تحصیلی ۱۳۹۹–۱۴۰۰ دانشگاه صنعتی امیرکبیر

به نام خدا

سوال ۱) پیاده سازی k-means با استفاده از

در این سوال می خواهیم با استفاده از داده های موجود برای این بخش و با استفاده از اسپارک، الگوریتم تکرار شونده یه k-means در این سوال می خواهیم با استفاده از پیاده سازی نماییم. فایل داده مورد نیاز برای این بخش حاوی f بطر است که هر سطر آن بیانگر سندی است که با استفاده از یک بر دار ویژگی f بعدی بازنمایی شده است. فایل های f و f نیز به ترتیب حاوی سنترویدهای اولیه f خوشه هستند که در f به صورت رندوم تعیین شده و در f این سنترویدها تا حد ممکن، با در نظر گرفتن معیار فاصله ی اقلیدسی، از یکدیگر فاصله دارند.

برای تمام قسمتهای این بخش، حداکثر تعداد تکرارها را برابر با ۲۰ و تعداد خوشهها را برابر با ۱۰ در نظر بگیرید.

الف) با در نظر گرفتن معی<mark>ار فاصله ی اقلیدسی</mark>، برای هر تکرار تابع هزینه را محاسبه نمایید. این عمل بدین معنی است که میبایست برای تکرار نظر گرفتن معیار فاصله ی اقلید الله و c1 و c1 استفاده کنید. الگوریتم k-means را روی دادههای فایل مله با استفاده از مقادیر فایلهای c1 و c2 اجرا کرده سپس مقادیر بدست آمده تابع هزینه را به صورت نمودار بر حسب تکرار از ۱ تا ۲۰ برای هر یک از دو مورد c1 و c2 رسم نمایید.

ب) درصد تغییر هزینه الگوریتم k-means بین اجرای صفرم و اجرای دهم را طبق عبارت $\frac{cost[0]-cost[10]}{cost[0]}$ ، با استفاده از مقادیر اولیه سنترویدها در cost[0] با یکدیگر مقایسه نمایید (معیار فاصله را در این سوال اقلیدسی در نظر بگیرید). توضیح دهید که کدام یک مقداردهی اولیه بهتری داشته است.

ج) مورد الف را این بار با در نظر گرفتن فاصله ی منهتن به عنوان تابع هزینه تکرار نمایید.

د) درصد تغییرات بین اجرای صفرم و دهم را این بار با در نظر گرفتن فاصلهی منهتن بدست آورده و نتایج را برای دو فایل c 1 و c2 مقاسه نمایید.

الگوریتم k-means تکرار شونده برای این تمرین در زیر ارائه شده است.

Algorithm 1 Iterative k-Means Algorithm

- 1: **procedure** Iterative k-Means
- 2: Select *k* points as initial centroids of the *k* clusters.
- 3: **for** iterations := 1 to MAX_ITER **do**
- 4: **for** cach point *p* in the dataset **do**
- 5: Assign point *p* to the cluster with closest centroid
- 6: end for
- 7: Calculate the cost for this iteration.
- 8: **for** each cluster *c* **do**
- 9: Recompute the centroid of c as the mean of all the data points assigned to c
- 10: **end for**
- 11: **end for**
- 12: end procedure

سوال ۲) Singular value decomposition

الف) اثبات کنید که مقادیر ویژه غیر صفر MM^T با مقادیر ویژه غیر صفر M^TM برابرند. همچنین بحث کنید که آیا بردارهای ویژه این دو نیز با یکدیگر برابرند یا خیر؟

 Σ و ماتریس M را به حاصلضرب $U\Sigma V^T$ میشکند که در آن دو ماتریس U و V، به صورت ستونی orthonormal بوده و Σ میشکند که در آن دو ماتریس Σ میشکند که در آن دو ماتریس میشکند که در آن دو ماتریس میشکند که در آن دو میشکند که در آن دو ماتریس میشکند که در آن دو میشکند که در آن در آ

ج) ماتریس زیر را در نظر بگیرید.

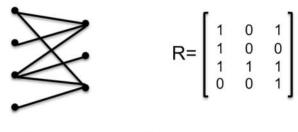
$$M = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 3 & 4 \\ 4 & 3 \end{bmatrix}$$

برنامهای بنویسید که SVD ماتریس M را محاسبه نماید. مقادیر بازگشتی برای ماتریسهای Σ و Σ را گزارش نمایید. سپس تجزیه ی مقدار ویژه حاصلضرب ماتریسی Σ Σ Σ Σ و امحاسبه کرده، مقادیر ویژه بدست آمده را به صورت نزولی مرتب کرده، لیست بردارهای ویژه را نیز مرتب سازی نمایید به صورتی که بردار ویژه متناظر با بزرگترین مقدار ویژه در ابتدای لیست باشد، لیست هر دو را گزارش دهید. در نهایت ارتباط میان ماتریس Σ و Σ بدست آمده با استفاده از تجزیه Σ ماتریس Σ و ماتریسهای بردارها و مقادیر ویژه گزارش شده در این بخش را بیان نمایید.

سوال ۳) سیستمهای توصیه گر

یک گراف دو بخشی کاربر-اقلام را در نظر بگیرید، در صورتی که بین کاربر U و قلم جنس I یک یال وجود داشته باشد، می گوییم کاربر I قلم جنس I را می پسندد. ماتریس نمرات برای این مجموعه از کاربرها و اقلام را با I نمایش می دهیم؛ هر سطر از I مربوط به یک کاربر بوده و هر ستون متعلق به یکی از اقلام است. در صورتی که کاربر i قلم جنس i را بپسندد I و در غیر این صورت I خواهد بود. فرض کنید I کاربر و I قلم جنس وجود دارد. حال ماتریس I را یک ماتریس قطری I تعریف می کنیم، I آمین عنصر قطری این ماتریس بیانگر تعداد عناصری است که کاربر I پسندیده است، به همین طریق ماتریس I را یک ماتریس قطری I تعریف می کنیم که I تعریف می کنیم که نامین عنصر قطری I تعداد کاربرانی هستند که قلم جنس I م را پسندیده اند. شکل زیر را به عنوان مثال مشاهده کنید.

Users Items



$$P = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad Q = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 3 \end{bmatrix}$$

الف) ماتریس شباهت کاربران را به صورت $T=R.R^T$ تعریف می کنیم. توضیح دهید T_{ij} و T_{ij} چه مفهومی دارد؟

ب) حال ماتریس شباهت اقلام را تعریف می کنیم، S_I یک ماتریس n^*n خواهد بود، عنصر موجود در سطر i و ستون j بیانگر شباهت i حال ماتریس شباهت اقلام i و i است. نشان دهید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i و است. نشان دهید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i و است. نشان دهید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i را می توان به فرم معادله (۱–۴) نوشت. توجه کنید i را می توان به فرم معادله (۱–۱) نوشت. توجه کنید i را می توان به نوش را می توان به

$$S_I = Q^{-1/2} R^T R Q^{-1/2} (1-f)$$

در ادامه همین سوال را برای ماتریس شباهت کاربران تکرار کنید، در اینجا عنصر سطر iام و ستون jام ماتریس S_U بیانگر شباهت کسینوسی میان کاربر iام و iام است. توجه کنید جواب نهایی باید یک عبارت بر اساس ماتریسهای i و iام است. توجه کنید جواب نهایی باید یک عبارت بر اساس ماتریسهای i و i باشد.

ج) روش توصیه بر اساس user-user collaborative filtering برای کاربر u می تواند به این صورت تعریف شود که برای همه اقلام مانند v معادله v معادله v شود و v قلم از اجناس که بیشترین v را داشته باشند به عنوان جواب بازگردانده شود.

$$r_{u.s} = \sum_{x \in users} \cos_s sim(x.u) \times R_{xs}$$
 (Y-4)

درس کلان دادهها

k (۳-۴) با استفاده از روش توصیه بر اساس item-item collaborative برای قلم جنس u با استفاده از معادله و استفاده از معادله قلم جنس جهت توصیه انتخاب می شود.

$$r_{u.s} = \sum_{x \in items} R_{xs} \times \cos_s sim(x.s)$$
 (7-4)

حال ماتریس توصیه به نام Γ را ابنگونه تعریف می کنیم که $arFi_{(i.j)}=r_{i.j}$ برای هر دو روش توصیه (اقلام-اقلام و کاربر-کاربر) ماتریس Γ را بر اساس Γ و Γ به دست آورید.

د) در این قسمت روشهای گفته شده را روی دادگان واقعی پیادهسازی می کنیم. دادگان مربوط به این قسمت شامل اطلاعات برنامههای تلوزیونی است، شامل ۹۹۸۵ کاربر و ۵۶۳ برنامه معروف. فایل user-shows.txt ماتریس R است که هر سطر آن نشانگر یک کاربر و هر ستون آن مرتبط با یک برنامه است، در صورتی که یک کاربر یک برنامه را بیشتر از سه ماه تماشا کند، درایه متناظر با برنامه و کاربر یک می شود. ستونها با space جدا شدهاند. فایل shows.txt شامل اسامی برنامهها به همان ترتیب ستونهای ماتریس R.

روشهای توصیه user-user collaborative filtering و user-user collaborative filtering برای کاربر ۱۵۰۰م از دادگان مقایسه خواهند شد، این کاربر را Alex مینامیم، توجه کنید ایندکس این کاربر در دادگان ۴۹۹ است. برای انجام این کار دادگان مقایسه خواهند شد، این کاربر را ۱۰۰۰ ورودی را برای سطر الکس صفر می کنیم، بدین معنی که دیگر نمی دانیم الکس کدام برنامهها را تماشا کرده است. بر اساس رفتار الکس با دیگر برنامهها ۱۰۰ پیشنهاد برای الکس پیدا می کنیم و سپس با برنامههای واقعی (که قبلا حذفشان کرده بودیم) مقایسه می کنیم تا میزان تطابق را به دست آوریم.

ا ماتریسهای P و Q را محاسبه کنید.

S مجموعه user-user collaborative filtering را برای Γ را برای user-user collaborative filtering محاسبه کنید. فرض کنید S مجموعه صد برنامه اول است، از تمام برنامههای موجود در S S برنامه که بیشترین شباهت با الکس را دارد پیدا کنید(نام برنامه را ذکر کنید و در صورتی که دو برنامه امتیاز یکسان داشتند، برنامه با ایندکس کمتر را انتخاب کنید).

۳- مانند S مجموعه تلوزیونی در مجموعه محاسبه کنید، از میان همه برنامههای تلوزیونی در مجموعه Γ مانند مانند ابرای الکس گزارش کنید. علاوه بر این امتیاز برنامههای انتخاب شده را نیز گزارش کنید.

توضيحات

- ❖ مهلت تحویل تمرین: ۱۴۰۰/۳/۲۶
- 💠 زمان ارائه: پس از تحویل مشخص خواهد شد.

نکاتی در مورد تحویل تمرین:

- ✓ فایل دادههای مورد نیاز تمرینها، به همراه فایل تمرین در سامانه بارگذاری شده است.
- ✓ خروجی کد ها و نتایج سوالات را درون گزارش بنویسید و از توضیح اضافی کد و موارد دیگر خودداری فرمایید (کد بدون گزارش ارزشی ندارد).
- ✓ فرمت تحویل: برای هر سوال یک پوشه ی جداگانه در نظر گرفته، کد و مواردی از قبیل خروجی برنامه و نمودارها را در آن ذخیره نمایید. این پوشه ها به همراه یک فایل report.pdf برای گزارش و توضیح سوالات، درون یک فایل فشرده شده با فرمت .zip
 ین بهتیبانی نمی شوند.
- ✓ در نظر داشته باشید کد های شما باید قابلیت اجرا در هنگام ارائه را داشته باشند. همچنین بر روی کد های خود کاملا مسلط باشید.
- ✔ توصیه میشود پیش از ارائه نیز مطالعهای روی کد و گزارش خود داشته باشید تا سوالات تدریس یاران را به راحتی پاسخ دهید.
- ✓ می توانید از گوگل برای رفع سوالات و مشکلات خود استفاده نمایید. در صورت رفع نشدن مشکل، می توانید سوالات خود را با
 تدریس یاران درس از طریق ایمیل زیر در میان بگذارید.
 - ❖ ایمیل تدریسیاران درس:

bdta00@gmail.com