



امتحان میان ترم درس داده کاوی

نیم سال دوم تحصیلی ۱۳۹۹-۰۰

دانشکده مهندسی کامپیوتر
مدرس: دکتر فضلی

نام: امیرحسین کارگران خوزانی شماره دانشجویی: ۹۹۲۰۱۱۱۹

سؤال ۱.

راه حل.

من راستش نمیدونستم هر ستونی دقیقا معنیش چیه و بیشتر اونهایی که میدونستم معنیش چیه و ممکنه کمک کنه را نگه داشتم. تعدادی از ستون ها را با ۰ پر کردم. و بعدش اونایی که مقدار NAN زیادی بیشتر از ۸۰۰ داشت را حذف کردم. سپس بقیه را با ۰ دوباره جایگذاری کردم. GOALS از مهمترین ستون ها می باشد چرا که برای مثال دروازه زبان ها یا مدافعان تعداد گل کمتری به ثمر رساندند و این برای مهاجم ها برعکس و مقدار زیادی است. یا مثلا CS فقط برای دروازه زبان است. کرنرها نیز موارد مهمی هستند که بازیکنان جلویی بیشتر آن را می زنن، همچنین سانت ها نیز. تقریبا بیشتر فیلدها مهم هستند و بیشتر مهم است که چگونه با آن ها رفتار شود و مقدار Nan آن ها با چه مقدار مناسبی پر شود.

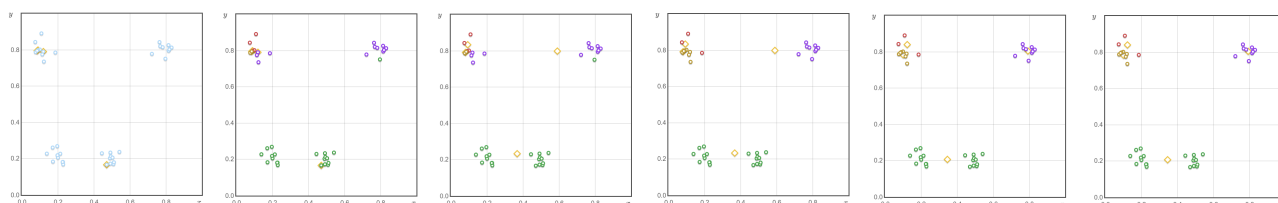
بخش ۵) ابتدا اثبات می کنیم که k-means به صورت یکنواخت $\psi(X^{(t)}) = \frac{1}{n} \sum_{j=1}^K \sum_{i=1}^n \|x_i, c_j\|^2$ را کاهش می دهد که $X^{(t)}$ تقسیم بندی فعلی $X_1^{(t)}, \dots, X_K^{(t)}$ با centroid های $c_1^{(t)}, \dots, c_K^{(t)}$ است. در هر مرحله از اجرای الگوریتم k-means، $\mathcal{A}(x_i)$ به صورت $\mathcal{A}(x_i) \leftarrow \arg \min_{j \in \{1, \dots, K\}} \|x_i - c_j\|^2$ تعریف می شود. که به $\mathcal{A}^{(t)}$ تابع انتساب گفته می شود. در این صورت از آنجا که $\mathcal{A}(x_i)$ مقدار $\|x_i - c_j\|^2$ را روی همه $j \in \{1, \dots, K\}$ کمینه می کند و $c_j^{(t+1)}$ نیز مقدار $\|x_i - c_j\|^2$ را روی همه $x_i \in X_j$ کمینه می کند، خواهیم داشت:

$$\begin{aligned} \psi(X^{(t)}) &\geq \sum_{j=1}^K \sum_{x_i \in X_j^{(t)}} \left\| x_i, c_{\mathcal{A}^{(t+1)}(x_i)}^{(t+1)} \right\|^2 \\ &\geq \sum_{j=1}^K \sum_{x_i \in X_j^{(t)}} \left\| x_i, c_j^{(t+1)} \right\|^2 \\ &\geq \psi(X^{(t+1)}) \end{aligned}$$

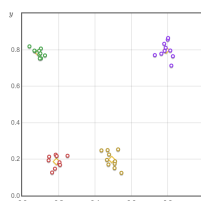
در این صورت اثبات می شود که این الگوریتم به طور یکنواخت مقدار $\frac{1}{n} \sum_{j=1}^K \sum_{i=1}^n \|x_i - c_j\|^2$ را کاهش می دهد. حال به اثبات صورت سوال گفته شده می پردازیم و تنها کافی است که نشان دهیم الگوریتم k-means پس از تعداد محدودی step متوقف می شود. از آنجا که تعداد محدودی تقسیم بندی وجود دارد و تعداد آن برابر $\binom{N}{k}$ است، در نتیجه دنباله $\psi(X^{(t)})_{t \in \mathbb{N}}$ تعداد محدودی مقدار خواهد داشت. این بدین معناست که حالتی وجود دارد که $\psi(X^{(t+1)}) = \psi(X^{(t)})$ باشد و در نتیجه در step شماره t داریم: $X^{(t+1)} = X^{(t)}$ در غیر این صورت برخی از element ها اشتباه دسته بندی شده است. در این اثبات صرفا یک باند نمایی برای مساله پیدا کردیم که نشان دهیم حتما در تعداد محدودی بار الگوریتم متوقف می شود.

بخش ۶) از آنجا که استفاده از منابع اینترنتی در طول امتحان مجاز است من از مثال این [link](#) استفاده کردم. هدف این لینک آن است که نشان دهد الگوریتم k-means در چه مثال هایی خوب عمل می کند و در چه مثال هایی بد. اما من مثالی که در آن این صفحه از آن به عنوان مثال خوب یاد کرده است صرفا استفاده کردم و نتیجه گیری ها کاملا مشاهدات و

استدلال‌های خودم است و تنها از اجرای آنالین ان استفاده کردم. همانگونه که در این صفحه گفته شده این مثالی است که در آن k-means به خوبی عمل می‌کند. اما اگر چندین بار الگوریتم را ران کنیم تا نقاطی که در ابتدا به عنوان نقاط تصادفی مرکز centroid ها initial می‌شود به گونه نامناسبی قرار گیرد و در نتیجه در ادامه اجرای الگوریتم در یک مینیمم محلی گیر می‌کند. برای مثال در این شکل ۱ ما انتظار داشتیم که ۴ کلاستر به شکل زیر ۲ در آید اما انتخاب نقاط اولیه باعث شده در یک مینیمم محلی گیر کنیم. و در نتیجه دو کلاستر پایینی با یکدیگر ادغام شده و یکی از کلاسترها به دو کلاستر تقسیم شود. دلیل این امر نیز آن است که الگوریتم روی همه تقسیم‌بندی‌های یک تابع محدب نیست و در نتیجه خروجی به مقداردهی اولیه بستگی دارد. به همین منظور توصیه می‌شود که چندین بار الگوریتم ران شود و خروجی مطلوب به عنوان جواب نهایی برداشته شود.



شکل ۱: اجرای الگوریتم کلاستریگ از سمت چپ به راست



شکل ۲: کلاستریگ مطلوب

سؤال ۲.

راه حل.

در نوتبوک به پیوست الگوریتم مقاله پیاده سازی شده است. در ادامه بحث می‌کنیم که الگوریتم این مقاله چگونه کار می‌دهد. اگر فرض کنیم که U_1, U_2, \dots به صورت iid بین ۰ و ۱ انتخاب شده باشند آنگاه اولین n ای که $S_n = \sum_{i=1}^n U_i > 1$ باشد را داشته باشیم. آنگاه $E(N)$ روی تمام n هایی که از این روش بدست آمده بگیریم برابر e خواهد بود. بنابراین یک تابع نوشتیم که n را برامون بدست بیاورد و در نتیجه تعداد زیادی از اون‌ها را با هم جمع کنه از طریق reduce و در نهایت هم تقسیم بر تعداد جمع این اعداد را بکنیم.

سؤال ۳.

راه حل.

حال اثبات می‌کنیم که چرا $E(N) = e$ می‌شود. ابتدا تابع احتمال متغیر N را تعیین می‌کنیم. به طور واضح N احتمال غیر صفر بر روی مجموعه $T = \{2, 3, 4, 5, \dots\}$ دارد و داریم: $\Pr(S_n > 1 \mid n \in T, \Pr(N = n)) = \Pr(S_{n-1} < 1) = \Pr(S_n < 1) - \Pr(S_n = 1)$

آخرین معادله به راحتی می‌تواند با نوشتن E_n برای رویداد $\{S_n < 1\}$ ($n = 1, 2, 3, \dots$) و با توجه به $\Pr(E_n^c \cap E_{n-1}) = \Pr(E_{n-1}) - \Pr(E_n)$ تایید شود.

حال نیاز است که مقدار $\Pr(S_n < 1)$ برای $n = 1, 2, 3, \dots$ محاسبه شود. این احتمال مقدار ۱ را برای $n = 1$ دارد. و اثبات می‌شود که $\Pr(S_n < 1) = 1/n!$ ($n = 1, 2, 3, \dots$) در حالیکه می‌توان گفت $\Pr(N = n) = (n-1)/n!$ ($n = 2, 3, 4, \dots$). بنابراین خواهیم داشت: $E(N) = \sum_{n=1}^{\infty} n \Pr(N = n) = e$.

قسمت اول) ابتدا به تعاریف می‌پردازیم و از روی تعاریف به نتیجه می‌رسیم. چه موقع یک FIS ماکسیمال است؟ اگر هیچ کدام از سوپرست های آن FIS دیگر بر تکرار نباشد و چه موقع یک FIS را بسته می‌نامیم؟ زمانی که هیچ سوپرستی از آن مقدار ساپورت برابر یا بزرگتری نداشته باشد. پس یک مجموعه بسته می‌تواند سوپرستی بر تکرار داشته باشد اما شرط بسته بودن نقض نشود به زبانی دیگر سوپرست پرتکراری می‌تواند داشته باشد ولی آن سوپرست مقدار ساپورت برابر یا بزرگتری نداشته باشد. پس FIS ماکسیمال زیر مجموعه FIS بسته است. همچنین اگر از طریق برهان خلف نیز به این مساله نگاه شود و در نظر بگیریم که FIS ماکسیمال زیر مجموعه FIS های بسته نیست پس باید شرط آن را نقض کند در صورتی که مجموعه FIS ماکسیمالی نمی‌توان یافت که سوپرست آن از آن مقدار ساپورت برابر یا بزرگتری داشته باشد چرا که با ماکسیمال بودن FIS در تناقض است.

قسمت دوم) همانگونه که این موضوع در کلاس نیز بحث شد اگر به فرمول confidence نگاه کنیم کاملاً متوجه آن خواهیم شد:

$$confidence(\{A \rightarrow B\}) = \frac{support(\{A \rightarrow B\})}{support(\{A\})} = \frac{freq(A, B)}{freq(A)} = \Pr(B|A)$$

اگر آیتی از A به B اضافه شود آنگاه صورت کسر تغییر نخواهد کرد. اما ممکن است باعث شود که مخرج کسر بزرگتر شود (چرا که یک آیت از مجموعه آن حذف شده است) و این موضوع باعث می‌شود مقدار confidence کاهش یابد و اگر تغییر نکند ثابت ماند که همان گزاره صورت سوال است.

قسمت سوم) بله از هر دو آن‌ها می‌توان استفاده کرد. با توجه به قسمت اول می‌توان گفت که اگر سمت چپ rule یک عبارت ماکسیمال پرتکرار داشته باشیم آنگاه دیگر نیازی به بررسی سوپرست‌های آن نیست چرا که دیگر پرتکرار نیستند تا در سمت چپ قرار گیرند. با توجه به قسمت دوم می‌توان rule هایی را که انتقال کالا باعث کم شدن confidence می‌شود را در صورتی که مقدار confidence آن‌ها از ترشولد کمتر می‌شود را در نظر نگرفت.

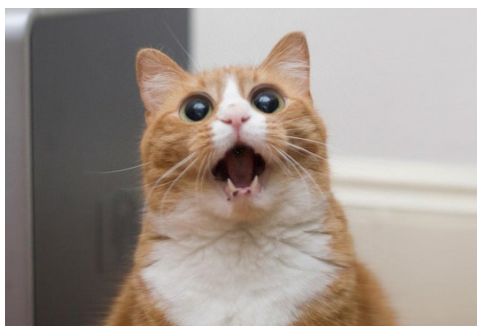
سؤال ۴.

راه حل.

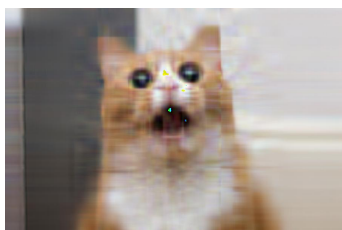
در قسمت ۱ با استفاده از کتابخانه pillow اقدام به بارگزاری عکس کردم و اندازه آن را چاپ کرده و کانال های مختلف رنگ را از یکدیگر نیز جدا کردم. در قسمت ۲ تابعی نوشتم که بتواند حاصلضرب اجزای مختلف svd را برای یک کانال حساب کند تا تصویر فشرده شده بدست آید. در قسمت ۳ برای بعد ۱۰۰ اقدام به محاسبه هر کانال کردم و در قسمت ۴ با استفاده از قسمت ۳ عکس کامل را ساختم این عکس را برای این پارامتر می‌توانید در شکل ۳ مشاهده کنید. این کار موجب

شد عکس از ۷۱ به ۳۸ کیلوبایت فشرده شود. سپس در قسمت ۵، با گام ۱۰ از ۱ تا ۲۰۰ حرکت کردم و عکس‌های مختلف را ساختم. همانگونه که در نوتبوک میبینید دو عکس مربوط به پارامتر ۵۱ و ۶۱ را گذاشتم که با یکدیگر تفاوت خاصی (چشمی) ندارند و هر دو از کیفیت مناسبی برخوردار هستند. بقیه عکس‌ها را می‌توانید در شکل ۴ مشاهده کنید. در قسمت ۶ حجم تک تک عکس‌ها را با کتابخانه os خواندم و نمودار مربوطه را کشیدم. همانجور که میبینید با اضافه شدن بعد حجم عکس نیز افزایش می‌یابد. دلیل آن نیز کم تر حذف شدن دادگان در تابع فشرده ساز است که پارامتر dim آن را کنترل می‌کند. در واقع می‌توان به جای این که مقدار $n * m * 3$ تعیین کننده سایز نسبی عکس باشد، مقدار $3 * (1 + m + n) * dim$ برابر سایز نسبی (بایت) عکس فشرده شده باشد. اولین عبارت حجم اصلی تصویر با ۳ کانال است که طول در عرض آن ضرب شده است. اما در دومی تنها برای ۳ کانال و ۳ ماتریس U, Z, V اعداد ذخیره می‌شوند که تعداد اعداد غیر صفر آن‌ها برابر عبارت دوم است. بدین صورت با افزایش dim این مقدار افزایش می‌یابد که در نمودار نیز مشخص است. اما از به جایی به می‌بینید که این نمودار شبیه کمتری پیدا می‌کند و از این فرمول حجم نسبی گفته شده کاملاً پیروی نمی‌کند این به آن دلیل است که خود عکس وقتی به صورت jpg ذخیره می‌شود توسط jpg از یک نوع فشرده‌ساز استفاده می‌کند و نمودار خود jpg نیز به شکل جذری است.

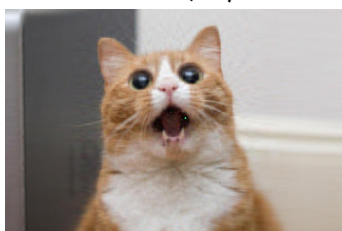
برای قسمت ۷، طبق اسلاید باید انرژی ماتریس singular-value در هر یک از حالت‌ها محاسبه شود و پارامتر بعد مناسب به نحوی انتخاب شود که حداقل ۹۰ درصد از انرژی را در خود داشته باشد. کد این قسمت نیز مانند بقیه قسمت‌ها در نوتبوک زده شده است و مقدار این فاصله انرژی با حد ۰/۰۰۱ مقدار انرژی محاسبه شده است. این مقدار dim برای هر یک از کانال‌ها به صورت مجزا حساب شده و سپس max آن در نظر گرفته شده است تا مقدار حد مورد نظر تضمین شود. همچنین راهکار دیگری که به ذهنم می‌رسد این است که ما دوست داریم حجم فایل‌ها کم شود در عین حال ارور بین عکس فشرده شده و عکس اصلی نیز کم باشد. با افزایش مقدار dim مشاهده می‌کنیم که این خطا کاهش می‌یابد اما در عوض حجم افزایش پیدا می‌کند (شکل خطا در بخش آخر نوتبوک کشیده شده است). در این صورت می‌توان یک تابع هدف تعریف کرد که با توجه به مقدار خطا و مقدار حجم فایل برای ما مقدار dim مناسب را پیدا کند. این تابع هدف دارای دو پارامتر مقدار خطا و حجم فایل در هر dim است که با استفاده از ۲ ضریب مثبت باید کمینه شود. هر کدام از ضریب‌ها در جه اهمیت مد نظر ما را برای هر پارامتر بیان می‌کند.



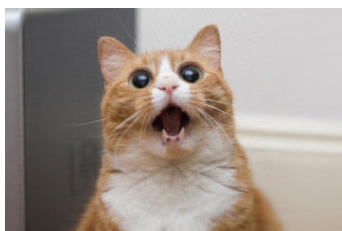
شکل ۳: بعد ۱۰۰



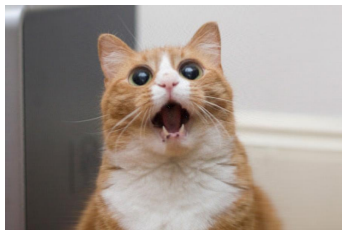
(ب) ابعاد ۱۱



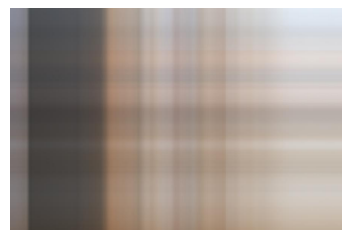
(د) ابعاد ۳۱



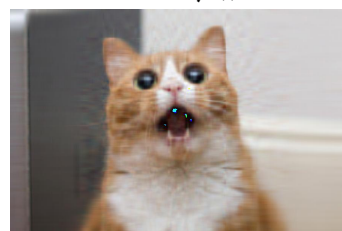
(و) ابعاد ۵۱



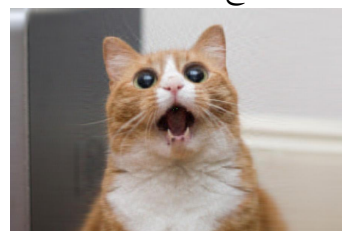
(ح) ابعاد ۷۱



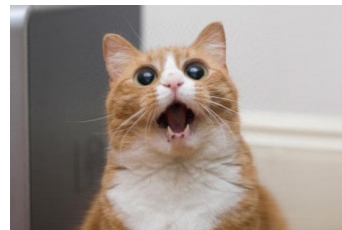
(آ) ابعاد ۱



(ج) ابعاد ۲۱



(ه) ابعاد ۴۱



(ز) ابعاد ۶۱

شکل ۴: شکل ها با پارامترهای مختلف بعد