

به نام خدا



دانشگاه شهید بهشتی
دانشکده مهندسی و علوم کامپیوتر

تحلیل مه داده‌ها
مجموعه تمرین سوم

استاد درس:

دکتر ملک

دانشجو:

پارسا جعفرقلی

۴۰۱۴۲۲۰۴۸

پاییز ۱۴۰۲

سوال اول (توجه کنید که تمام راه حل مرحله به مرحله آورده شود) (۲۰ نمره)

دوازده سبد شامل سه آیتم وجود دارد:

$\{1,2,3\}$ $\{2,3,4\}$ $\{3,4,5\}$ $\{4,5,6\}$ $\{1,3,5\}$ $\{2,4,6\}$

$\{1,3,4\}$ $\{2,4,5\}$ $\{3,5,6\}$ $\{1,2,4\}$ $\{2,3,5\}$ $\{3,4,6\}$

در صورتی که support برابر ۴ باشد و تابع hash برابر $(i*j)\%9$ باشد:

الف) با استفاده از الگوریتم PCY کدام جفت‌ها در Pass2 شمارش می‌شوند.

ب) این بار از الگوریتم multistage برای حل این مسأله استفاده کنید. تابع hash دوم برابر $(i+j)\%11$ است. کدام جفت‌ها در pass3 شمارش می‌شوند؟

ج) آیا انتخاب تابع hash $(i+j)\%11$ به عنوان دومین تابع hash انتخاب مناسبی است؟

Freq ۱: $\{۱\}$, $\{۲\}$, $\{۳\}$, $\{۴\}$, $\{۵\}$, $\{۶\}$

$h(۱,۲) = ۲$ $h(۱,۳) = ۳$ $h(۱,۴) = ۴$ $h(۱,۵) = ۵$ $h(۱,۶) = ۶$ $h(۲,۳) = ۶$

$h(۲,۴) = ۸$ $h(۲,۵) = ۱$ $h(۲,۶) = ۳$ $h(۳,۴) = ۳$ $h(۳,۵) = ۶$ $h(۳,۶) = ۰$

$h(۴,۵) = ۲$ $h(۴,۶) = ۵$ $h(۵,۶) = ۳$

الف) $h(۱,۳)$ و $h(۲,۶)$ و $h(۳,۴)$ و $h(۵,۶)$ وارد قسمت pass۲ می‌شوند.

$$h(5,6) = 0 \qquad h(3,4) = 7 \qquad h(2,6) = 8 \qquad h(1,3) = 4$$

ب) هیچ کدام داخل یک باکت نیوفتادند.

ج) ارزیابی تابع hash دوم

- تابع hash خوب تابعی است که توزیع یکنواخت داشته باشد تا collision کمتری داشته باشیم.

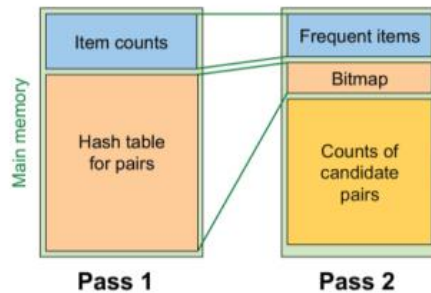
- ($h(i, j) = (i + j) \% 11$) ممکن است خوب نباشد اگر توزیع یکنواخت نداشته باشیم، به خصوص وقتی که فضای ممکن برای سطل‌ها به اندازه مقدار مد نظر نباشد.

- باید توزیع جفت‌ها در این تابع hash مورد آزمایش و بررسی قرار گیرد تا اینکه بتوان نتیجه گرفت آیا انتخاب مناسبی است یا خیر.

بدیهی است که برای رسیدن به پاسخ‌های دقیق‌تر، نیازمند اجرا کردن فرآیندهای توضیح داده شده بر روی ماشین حساب یا کدنویسی برای شمارش دقیق و تعیین bucketهای hash فرکوئنت هستیم. این کار نیازمند جمع‌آوری و تجزیه و تحلیل داده‌های خام است که در این متن امکان‌پذیر نیست.

سوال دوم (برای گزینه درست راه حل آورده شود. اعداد صحیح با ۴ بایت نمایش داده می‌شوند.) (۱۰ نمره)

در الگوریتم PCY اگر فضای مورد نیاز برای شمارش آیتم‌ها 4 MB باشد و تعداد باکت‌ها برابر $8 * 10^9$ باشد، به طور تقریبی چه فضایی از Main Memory برای شمارش جفت آیتم‌ها در Pass 2 قرار می‌گیرد (قسمت زرد رنگ در شکل)؟



31 GB (د)

28 GB (ج)

32 GB (ب)

الف) اطلاعات مسأله ناقص است

هر کدام از هش تیبل‌ها برای محاسبه به یک فضای ۴ بایتی نیاز دارند. پس فضایی که کل هش تیبل می‌گیرد برابر $4 * 8 * 10^9$ است که ۳۲ گیگابایت می‌باشد. از طرفی فضای بردار Bitmap یک فضای $1/32$ از فضای هش تیبل را پر می‌کنند بنابراین یک بردار ۱ گیگابایتی است. با توجه به این که اندازه شمارش ایت‌م ۴ مگابایت است و فرکونت ایت‌م نیز کوچک تر از آن است (تنها مقادیر فرکونت را بر می‌گرداند) برای همین زیاد در محاسبات اهمیت ندارند. پس مقدار باقی مانده می‌شود هش مپ قسمت اول منهای بیت مپ قسمت دوم که می‌شود ۳۱ گیگابایت.

سوال سوم (۲۰ نمره)

مجموعه داده‌ای را در فضای دو بعدی با توجه به نقاط $A(1,2)$, $B(3,4)$, $C(5,6)$, $D(7,8)$, $E(9,10)$ در نظر بگیرید. الگوریتم k -means را با فرض ۲ خوشه ($k=2$) و نقاط ابتدایی $A(1,2)$ و $D(7,8)$ برای یک بار روی مجموعه اعمال کرده و مقادیر مراکز خوشه‌ها را به روزرسانی کنید. در نهایت مراکز خوشه جدید و نقاط اختصاص داده شده به هر کلاستر را نمایش دهید.

بیاپید الگوریتم k -means را برای یک تکرار با $k=2$ خوشه و نقاط شروع $A(1,2)$ و $D(7,8)$ پیاده کنیم. این مراحل به صورت زیر است:

مرحله ۱: تعیین مراکز اولیه خوشه

خوشه ۱: $C_1 = A(1,2)$

خوشه ۲: $C_2 = D(7,8)$

مرحله ۲: اختصاص دادن نقاط به نزدیکترین مرکز خوشه

نقاط B , C , و E را در نظر گرفته و فاصله هر کدام را از دو مرکز خوشه اندازه‌گیری می‌کنیم. برای فاصله از فاصله اقلیدسی استفاده می‌کنیم:

فاصله بین دو نقطه (x_1, y_1) و (x_2, y_2) از رابطه زیر محاسبه می‌شود:

$$\text{فاصله} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

مرحله ۳: محاسبه فاصله‌ها و اختصاص به خوشه‌ها

با استفاده از رابطه بالا، فاصله هر نقطه (B, C, E) را از هر دو مرکز خوشه (C۱ و C۲) محاسبه و مقایسه می‌کنیم:

$$\text{فاصله B از C۱} : \sqrt{(3-1)^2 + (4-2)^2} = \sqrt{8}$$

$$\text{فاصله B از C۲} : \sqrt{(7-3)^2 + (8-4)^2} = \sqrt{32}$$

B به C۱ نزدیک‌تر است.

$$\text{فاصله C از C۱} : \sqrt{(5-1)^2 + (6-2)^2} = \sqrt{32}$$

$$\text{فاصله C از C۲} : \sqrt{(7-5)^2 + (8-6)^2} = \sqrt{8}$$

C به C۲ نزدیک‌تر است.

$$\text{فاصله E از C۱} : \sqrt{(9-1)^2 + (10-2)^2} = \sqrt{128}$$

$$\text{فاصله E از C۲} : \sqrt{(7-9)^2 + (8-10)^2} = \sqrt{8}$$

E به C۲ نزدیک‌تر است.

جدول اختصاص نقاط به خوشه‌ها:

خوشه ۱ فاصله از C۱ خوشه ۲ فاصله از C۲

$$A(۱,۲) \cdot$$

$$B(۳,۴) \sqrt{۸}$$

$$C(۵,۶) \sqrt{۳۲}$$

$$D(۷,۸) \cdot$$

$$E(۹,۱۰) \sqrt{۸}$$

با توجه به جدول بالا، B به C۱، در حالی که C و E به C۲ اختصاص داده شده‌اند.

گام چهارم: به روزرسانی مراکز خوشه

میانگین مختصات نقاط اختصاص داده شده به هر خوشه محاسبه می‌شود تا مختصات مراکز جدید خوشه‌ها تعیین شوند.

برای خوشه ۱ که شامل نقاط A(۱,۲) و B(۳,۴) می‌شود:

$$\text{میانگین مختصات } x: ۲ = ۲/(۱+۳)$$

$$\text{میانگین مختصات } y: ۳ = ۲/(۲+۴)$$

پس، مرکز جدید خوشه ۱، C۱_جدید { می‌شود: (۳,۲)

برای خوشه ۲ که شامل نقاط $C(۵,۶)$ ، $D(۷,۸)$ و $E(۹,۱۰)$ می‌شود:

$$۷ = ۳/(۵+۷+۹) :x \text{ میانگین مختصات}$$

$$۸ = ۳/(۶+۸+۱۰) :y \text{ میانگین مختصات}$$

پس، مرکز جدید خوشه ۲، $C_2\{\text{جدید}\}$ می‌شود: $(۸,۷)$

در نتیجه، مراکز جدید خوشه‌ها به شکل زیر است:

$$(۲,۳) = C_1\{\text{جدید}\} : ۱ \text{ - مرکز جدید خوشه}$$

$$(۷,۸) = C_2\{\text{جدید}\} : ۲ \text{ - مرکز جدید خوشه}$$

سوال چهارم (۵ نمره)

فرض کنید می‌خواهیم الگوریتم Balance را برای تخصیص تبلیغات به دو تبلیغ کننده A و B اجرا کنیم. تبلیغ کننده A روی کوثری‌های x و y قیمت می‌دهد و تبلیغ دهنده B روی کوثری‌های x و z. هر دو نیز بودجه‌ای برابر ۲ دلار دارند. آیا الگوریتم می‌تواند برای کوثری zxyy تخصیص بهینه را انجام دهد؟

برای حل این مثال با استفاده از الگوریتم Balance، ما باید نحوه توزیع تبلیغات بر روی کوثری‌ها را بررسی کنیم. ما بودجه مساوی برای تبلیغ‌کننده‌های A و B داریم (هر کدام ۲ دلار) و می‌خواهیم برای سلسله‌ای از کوثری‌ها تبلیغ‌ها را تخصیص دهیم که به ترتیب zxyy هستند. برای این کار، هر تبلیغ‌دهنده به طور عادلانه باید چهار تبلیغ از سری کوثری‌ها را تقسیم کند. با این فرض که هر کلیک دقیقاً ۱ دلار هزینه داشته باشد.

نحوه اجرا به این صورت خواهد بود:

۱. کوثری z فقط توسط B پوشش داده می‌شود، بنابراین اولین تبلیغ z به B تخصیص داده می‌شود. حالا B دقیقاً ۱ دلار از بودجه‌اش را خرج کرده است.

۲. کوثری x هم توسط A و هم توسط B پوشش داده می‌شود. از آنجا که ما می‌خواهیم بودجه را متعادل نگه داریم، این کوثری باید به A داده شود چون B در حال حاضر ۱ دلار هزینه کرده و A هنوز هزینه‌ای نداشته است. حالا A و B هر کدام ۱ دلار از بودجه خود را خرج کرده‌اند.

۳. دو کوثری y بعدی فقط توسط A پوشش داده می‌شود، بنابراین هر دو را به A تخصیص می‌دهیم. حالا A کل بودجه ۲ دلاری خود را خرج کرده است.

نتیجه تخصیص بهینه به شکل زیر خواهد بود:

- کوثری z - تبلیغ کننده B (۱ دلار باقی مانده برای B)

- کوثری x - تبلیغ کننده A (۱ دلار باقی مانده برای A)

- کوثری y - تبلیغ کننده A (بودجه A تمام شد)

- کوثری y - تبلیغ کننده A (بودجه A تمام شد)

در این نقطه، تبلیغ کننده A دیگر بودجه‌ای ندارد، ولی هنوز یک کوثری x باقی مانده است که هر دو تبلیغ کننده می‌توانند پاسخ دهند و ما همچنان ۱ دلار بودجه از تبلیغ کننده B باقی مانده داریم که می‌توانیم خرج کنیم. پس، آخرین کوثری x به تبلیغ کننده B داده می‌شود. این تبلیغ باعث می‌شود که بودجه‌ی B هم تمام شود.

نتیجه نهایی تخصیص بودجه تبلیغات به شکل زیر خواهد بود:

- تبلیغ کننده A : کوثری x و دو بار کوثری y تخصیص داده شد (بودجه صرف شده = ۲ دلار)

- تبلیغ کننده B : کوثری z و کوثری x تخصیص داده شد (بودجه صرف شده = ۲ دلار)

بنابراین، الگوریتم Balance بهینه این A, A, B, B است که تولید نشد.

سوال پنجم (۲۰ نمره)

یک ناشر آگهی، سه آگهی را برای قرار دادن در هر صفحه به ترتیب از بالا انتخاب می‌کند. CTR در هر موقعیت برای هر تبلیغ کننده متفاوت است و هر تبلیغ کننده برای هر موقعیت CTR متفاوتی دارد. هر تبلیغ کننده برای تعداد کلیک‌ها پیشنهاد قیمت می‌دهد و یک بودجه روزانه دارد. هنگامی که یک کلیک رخ می‌دهد، تبلیغ کننده مبلغی را که پیشنهاد داده است پرداخت می‌کند. در زیر جدولی از پیشنهادها، CTR برای موقعیت‌های ۱، ۲، ۳ و بودجه برای هر تبلیغ کننده آمده است.

Advertiser	Bid	CTR1	CTR2	CTR3	Budget
A	\$.10	.015	.010	.005	\$1
B	\$.09	.016	.012	.006	\$2
C	\$.08	.017	.014	.007	\$3
D	\$.07	.018	.015	.008	\$4
E	\$.06	.019	.016	.010	\$5

برای پیدا کردن تعداد کلیک‌های هر تبلیغ دهنده باید میزان درآمد مورد انتظار برای هر موقعیت از هر تبلیغ کننده محاسبه کنیم و موقعیت‌های تبلیغاتی را بر اساس بالاترین درآمد مورد انتظار تخصیص دهیم. فرمول انتظار میزان درآمد برای هر موقعیت به صورت CTR مربوط به آن موقعیت ضربدر پیشنهاد قیمت (bid) برای هر کلیک است.

ابتدا برای هر تبلیغ کننده، میزان درآمد مورد انتظار برای هر موقعیت را محاسبه کنیم:

قرار است ناشر آگهی‌ها را بر اساس بالاترین میزان درآمد مورد انتظار انتخاب کند و هر تبلیغ کننده فقط یکبار می‌تواند برای یک موقعیت انتخاب شود.

۱. موقعیت ۱:

$$A: \$0.10 * 0.15 = \$0.015 -$$

$$B: \$0.09 * 0.16 = \$0.0144 -$$

$$C: \$0.08 * 0.17 = \$0.0136 -$$

$$D: \$0.07 * 0.18 = \$0.0126 -$$

$$E: \$0.06 * 0.19 = \$0.0114 -$$

تبلیغ کننده A بالاترین میزان درآمد مورد انتظار را برای موقعیت ۱ دارد.

۲. موقعیت ۲:

$$B: \$0.09 * 0.12 = \$0.0108 -$$

$$C: \$0.08 * 0.14 = \$0.0112 -$$

$$D: \$0.07 * 0.15 = \$0.0105 -$$

$$E: \$0.06 * 0.16 = \$0.0096 - \text{ (A را حذف می‌کنیم چرا که برای موقعیت ۱ انتخاب شده)}$$

تبلیغ کننده C بالاترین میزان درآمد مورد انتظار را برای موقعیت ۲ دارد.

۳. موقعیت ۳:

$$B: \$0.09 * 0.006 = \$0.00054 -$$

$$D: \$0.07 * 0.008 = \$0.00056 -$$

$$E: \$0.06 * 0.010 = \$0.0006 - \text{ (A و C را حذف می‌کنیم چون برای موقعیت‌های قبل انتخاب شده‌اند)}$$

تبلیغ‌کننده E بالاترین میزان درآمد مورد انتظار را برای موقعیت ۳ دارد.

حالا، باید تعداد کلیک‌هایی که هر تبلیغ دهنده در یک روز بدست می‌آورد محاسبه کنیم. هر تبلیغ‌دهنده برای تعداد کلیک‌های بدست آمده مبلغ پیشنهادی خود را پرداخت می‌کند تا زمانی که بودجه‌شان تمام شود.

برای مثال، اگر فرض کنیم که در یک روز ۱۰۰ کلیک وجود دارد:

- تبلیغ‌کننده A برای هر کلیک \$۰.۱۰ پرداخت می‌کند و بودجه‌اش \$۱ است، پس تبلیغ‌کننده A می‌تواند ۱۰ کلیک (۱ / ۰.۱۰) بدست آورد.

- تبلیغ‌کننده C برای هر کلیک \$۰.۰۸ پرداخت می‌کند و بودجه‌اش \$۳ است، پس تبلیغ‌کننده C می‌تواند ۳۷.۵ کلیک (۳ / ۰.۰۸) بدست آورد، اما چون کلیک‌ها نمی‌توانند نیمه باشند، ۳۷ کلیک تمام کلیک‌هایی خواهد بود که تبلیغ‌کننده C می‌تواند بدست آورد.

- تبلیغ‌کننده E برای هر کلیک \$۰.۰۶ پرداخت می‌کند و بودجه‌اش \$۵ است، پس تبلیغ‌کننده E می‌تواند ۸۳.۳۳ کلیک (۵ / ۰.۰۶) بدست آورد، اما باز هم باید به تعداد تمام کلیک‌ها گرد شود پس ۸۳ کلیک بدست خواهد آمد.