

به نام خدا  
دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)  
دانشکده مهندسی کامپیوتر



دانشگاه صنعتی امیرکبیر  
(پلی تکنیک تهران)

## گزارش تمرین اول درس کلان داده‌ها

استاد درس: دکتر مصطفی حقیر چهرقانی

دانشجو: فاطمه غلامزاده

۹۹۱۳۱۰۰۳

نیم سال دوم ۱۳۹۹-۱۴۰۰

درس کلاس داده‌ها - تمرین اول - فاطمه غلام‌نژاد - ۹۶۱۳۱۰۰۳

سوال ۱ (الف - ۱) چه مقدار association rule متفاوت

(حتی با support صفر) می‌توان از این سبدهای خرید استخراج کرد؟

پاسخ:

تعداد آئتم‌های متفاوت ۴ عدد است. تعداد افزازهای یک مجموعه ۴ عضوی برای ۵ آیتم

ضربدر ۲ می‌کنیم تا مقدار rule‌های متفاوت بدست آید. برای راحتی در محاسبات آئتم‌های

ممانع را با اعداد ۱ تا ۴ نشان می‌دهیم:

حالت اول: قانون‌هایی به شکل  $1 \rightarrow 2$  یا  $3 \rightarrow 5$ :

$$\binom{4}{1} \times 2 = 15 \times 2 = 30$$

حالت دوم: قانون‌هایی به شکل  $1 \rightarrow 2, 3$

$$\binom{4}{2} \times \binom{3}{1} \times 2 = 20 \times 3 \times 2 = 120$$

حالت سوم: قانون‌هایی به شکل  $1 \rightarrow 2, 3, 4$

$$\binom{4}{3} \times \binom{3}{1} \times 2 = 15 \times 4 \times 2 = 120$$

حالت چهارم: قانون‌هایی به شکل  $1, 2 \rightarrow 3, 4$

$$\binom{4}{2} \times \binom{3}{1} = 15 \times 4 = 60$$

حالت پنجم: قانون‌هایی به شکل  $1 \rightarrow 2, 3, 4, 5$

$$\binom{4}{4} \times \binom{3}{1} \times 2 = 4 \times 5 \times 2 = 40$$

حالت ششم:  $1, 2 \rightarrow 3, 4, 5$

$$\binom{4}{5} \times \binom{3}{2} \times 2 = 4 \times 10 \times 2 = 80$$

حالت هفتم:  $1 \rightarrow 2, 3, 4, 5, 6$

$$\binom{4}{6} \times \binom{3}{3} \times 2 = 1 \times 4 \times 2 = 8$$

۱, ۲ → ۳, ۴, ۵, ۶

• حالت هشتم:

$$\binom{4}{4} \binom{4}{2} \times 2 = 15 \times 2 = 30$$

۱, ۲, ۳ → ۴, ۵, ۶

• حالت نهم:

$$\binom{4}{4} \binom{4}{3} = 1 \times 4 = 4$$

در حالت هایی که تعداد اعضای آنتیم ست ها در دو طرف قانون یکسان است

منبر ۲ نداریم چون حالت تکراری ایجاد می شود.

مجموع ۹ حالت:

$$30 + 120 + 120 + 90 + 40 + 120 + 12 + 30 + 4$$

$$= 402$$

← مقدار association rule های متفاوت

الف- ۲) بیشترین اندازه ی آنتیم ستی که می توانیم استخراج کنیم؟

با توجه به جدول داده شده، آنتیم ست های شماره ۴ و ۹ دارای بیشترین اندازه هستند

که ۴ است ← پس بیشترین اندازه ی Itemset = 4

الف- ۳) عبارتی برای بیشینه مقدار Itemset های با اندازه ۳ که می توانیم استخراج کنیم؟

$$\binom{4}{3} = 4$$

← انتخاب ۳ عضو از ۴ آنتیم متفاوت

الف... ۴) نسبت به الگوریتم apriori عمل می‌کنیم. ابتدا ساپورت Itemset ها

تک عضوی را می‌یابیم:  $sup(شیر) = 4$

$sup(کره) = 5$   $sup(نوشابه) = 4$

$sup(سیکویست) = 4$   $sup(دستمال) = 7$

$sup(نان) = 4$

عدد به سبب Itemset های ۲ عضوی داریم، سعی می‌کنیم از آن‌ها شروع کنیم که ترکیب

Itemset های تک عضوی با ساپورت بالا هستند:

$sup(شیر، دستمال) = 5$

با توجه به اینکه ساپورت {شیر، دستمال}

عدد ۵ بدست آمده می‌توانیم آئیم ست‌های ۲ تایی که یکی از اعضای آن ساپورت کمتر از ۵ دارند

را به طور کلی حذف کنیم و شمارش نکنیم چون قطعاً ساپورت آن‌ها از ۵ کمتر خواهد شد.

پس ترکیب‌های ۲ تایی که سیکویست، نان و نوشابه را در خود دارند رد نخواهیم کرد.

$sup(کره و شیر) = 3$

$sup(دستمال، کره) = 3$

پس Itemset با اندازه ۲ که بیشترین ساپورت را دارند همان {شیر، دستمال} است.

آئیم ست‌های با اندازه ۲ بزرگتر از ۲ هم قطعاً ساپورت کمتری از {شیر، دستمال} دارند چون

فقط یک کاندیدای ۳ عضوی داریم: {شیر، کره، دستمال} که ساپورت آن هم ۳ است.

پایان کجایی: {شیر، دستمال}



$$\text{Conf}(A \rightarrow B) = \frac{\text{sup}(A \cup B)}{\text{sup}(A)} \quad \text{الف- ۵}$$

$$\text{Conf}(B \rightarrow A) = \frac{\text{sup}(B \cup A)}{\text{sup}(B)}$$

$$\text{Conf}(A \rightarrow B) = \text{Conf}(B \rightarrow A) \Rightarrow \frac{\text{sup}(A \cup B)}{\text{sup}(A)} = \frac{\text{sup}(A \cup B)}{\text{sup}(B)}$$

$$\Rightarrow \text{sup}(A) = \text{sup}(B)$$

پس کافی است سایزِ اِنتِگِرتِستِ های A و B با هم برابر باشد. مثال:

بیسکوئیت  $\rightarrow$  نان

$$\text{Conf}(\text{نان} \rightarrow \text{بیسکوئیت}) = \frac{\text{sup}(\{\text{نان}, \text{بیسکوئیت}\})}{\text{sup}(\text{نان})} = \frac{1}{4} \quad \checkmark$$

$$\text{Conf}(\text{بیسکوئیت} \rightarrow \text{نان}) = \frac{\text{sup}(\{\text{نان}, \text{بیسکوئیت}\})}{\text{sup}(\text{بیسکوئیت})} = \frac{1}{4} \quad \checkmark$$

ب- ۱) خیر معیار confidence مقارن نیست. مثل نقق:

|      |   |
|------|---|
| D    | ۱ |
| C    | ۲ |
| A, B | ۳ |
| A    | ۴ |
| A, B | ۵ |

$$\text{Conf}(A \rightarrow B) = \frac{\text{sup}(A \cup B)}{\text{sup}(A)} = \frac{2}{3}$$

$$\text{Conf}(B \rightarrow A) = \frac{\text{sup}(A \cup B)}{\text{sup}(B)} = \frac{2}{2} = 1$$

$$\frac{2}{3} \neq 1$$

ب-۲) به Lift مقارن است. مقدار کل ترانسجسها را  $N$  در نظر می گیریم:

اثبات:

$$\text{Lift}(A \rightarrow B) = \frac{\text{Conf}(A \rightarrow B)}{S(B)}$$

$$= \frac{\frac{\text{sup}(A \cup B)}{\text{sup}(A)}}{\frac{\text{sup}(B)}{N}} = \frac{\text{sup}(A \cup B) \times N}{\text{sup}(A) \times \text{sup}(B)}$$

$$= \frac{\frac{\text{sup}(A \cup B)}{\text{sup}(B)}}{\frac{\text{sup}(A)}{N}} = \frac{\text{Conf}(B \rightarrow A)}{S(A)} = \text{Lift}(B \rightarrow A)$$

ب-۳) Conviction مقارن نیست. مثال نقض:

|      |   |
|------|---|
| D, B | ۱ |
| A, B | ۲ |
| A, B | ۳ |
| A    | ۴ |
| B    | ۵ |

$$\text{Conf}(A \rightarrow B) = \frac{۲}{۳}$$

$$\text{Conv}(A \rightarrow B) = \frac{1 - ۴/۵}{1 - ۲/۳} = \frac{۳}{۵}$$

$$\text{Conf}(B \rightarrow A) = \frac{۲}{۴}$$

$$\text{Conv}(B \rightarrow A) = \frac{1 - ۳/۵}{1 - ۲/۴} = \frac{۴}{۵}$$

$$\frac{۳}{۵} \neq \frac{۴}{۵} \Rightarrow \text{Conv}(A \rightarrow B) \neq \text{Conv}(B \rightarrow A)$$

## سوال ۲) کاربرد Association rules

**الف)** ۵ قانون ابتدایی به صورت زیر است که در آن مجموعه  $[A', B']$  به معنای قانون  $A \rightarrow B$  می باشد و عدد نوشته شده در روبه روی هر قانون، میزان confidence برای آن قانون می باشد.

```
[[['DAI93865', 'FRO40251'], 1.0],  
[['GRO85051', 'FRO40251'], 0.999176276771005],  
[['GRO38636', 'FRO40251'], 0.9906542056074766],  
[['ELE12951', 'FRO40251'], 0.9905660377358491],  
[['DAI88079', 'FRO40251'], 0.9867256637168141]]
```

**ب)** ۵ قانون ابتدایی به صورت زیر است که در آن مجموعه  $[A', B', C']$  به معنای قانون  $A, B \rightarrow C$  می باشد و عدد نوشته شده در روبه روی هر قانون، میزان confidence برای آن قانون می باشد.

```
[[['DAI23334', 'ELE92920', 'DAI62779'], 1.0],  
[['DAI31081', 'GRO85051', 'FRO40251'], 1.0],  
[['DAI55911', 'GRO85051', 'FRO40251'], 1.0],  
[['DAI62779', 'DAI88079', 'FRO40251'], 1.0],  
[['DAI75645', 'GRO85051', 'FRO40251'], 1.0]]
```

## سوال ۳) LSH

**الف)**

میانگین زمان جستجو در حالت خطی :

**Average LSH Time: 0.997 Seconds**

میانگین زمان جستجو با استفاده از LSH :

**Average Linear Search Time: 14.204 Seconds**

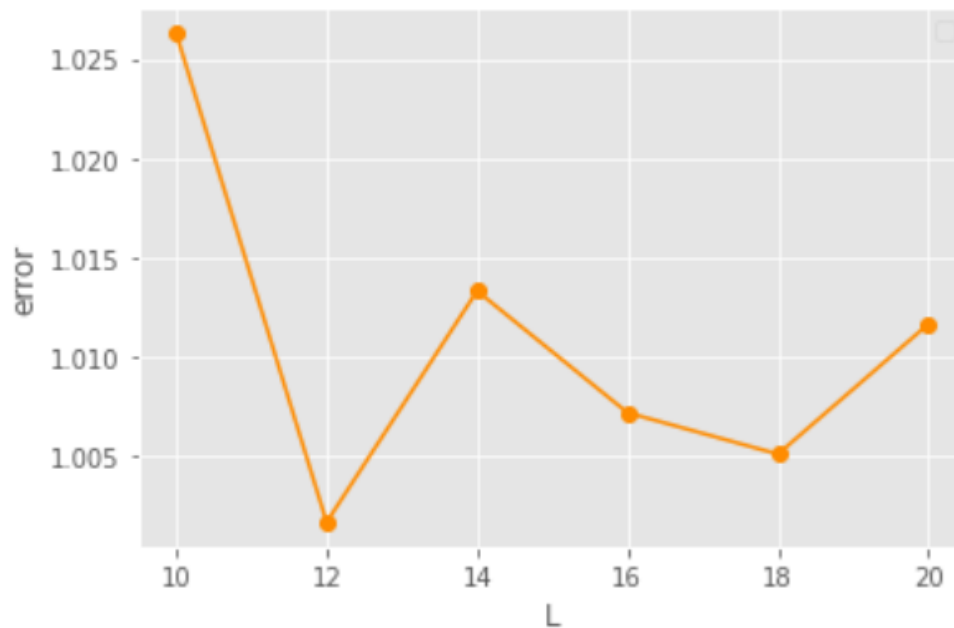
**مقایسه :** همانطور که مشاهده می شود میانگین زمان جستجو با استفاده از LSH نسبت به جستجوی خطی کاهش نسبتاً زیادی داشته است.

ب) مقدار خطای به دست آمده:

**Error: 1.023511**

ج)

نمودار مقدار خطا به صورت تابعی از L:



مقادیر خطا:

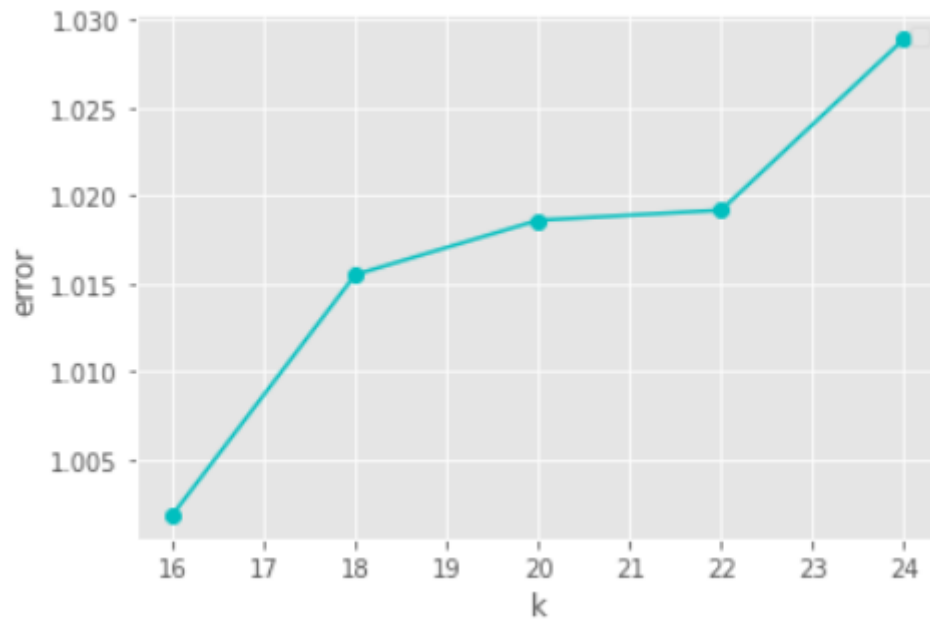
[1.0262867364053494, 1.0016595368482275,  
1.0133086954891755, 1.0071734390323535,  
1.0050886944625834, 1.0116438905081788]

تحلیل نمودار:

با افزایش مقادیر L مقدار خطا به طور کلی سیر نزولی دارد زیرا تعداد توابع hash افزایش می یابد.



نمودار مقدار خطا به صورت تابعی از  $k$ :



مقادیر خطا:

```
[1.0018693720851233, 1.01552212270779,  
1.018614931861614, 1.0191862558990539,  
1.0289129370499819]
```

تحلیل نمودار:

با افزایش  $k$  مقدار خطا افزایش یافته است چون در این حالت یافتن همسایه‌ها توسط LSH سخت‌تر می‌شود.

(د)



تصویر مورد سوال :

همسایه‌های به دست آمده با جستجوی خطی :



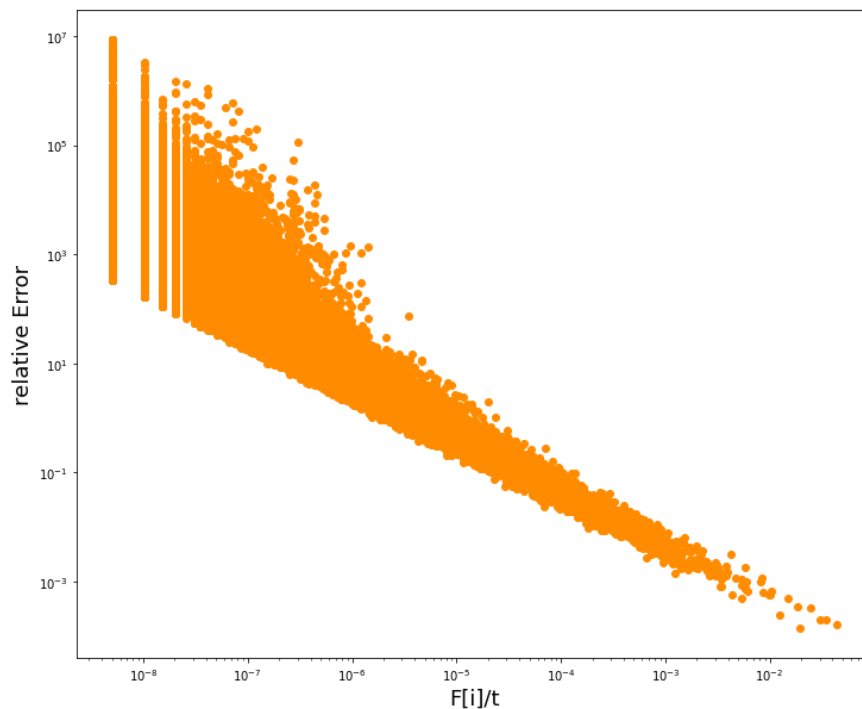
همسایه‌های به دست آمده با جستجوی LSH :



**مقایسه :** با توجه به جستجوی انجام شده و همسایه‌های یافت شده، با بررسی ایندکس همسایه‌ها مشخص گردید که در ۸ مورد LSH همسایه‌های یکسانی با linear پیدا کرده است و فقط در ۲ مورد تفاوت دارد. با توجه به تصاویر مشاهده می‌شود که آن ۲ مورد هم بسیار به هم شبیه هستند و همسایه نزدیکی برای تصویر مورد سوال توسط LSH گزارش شده است. بنابراین می‌توان نتیجه گرفت با توجه به نتایج خوبی که LSH دارد و زمان جستجوی بسیار کمتر آن نسبت به روش خطی، یک روش مناسب و به صرفه برای یافتن همسایه‌های نزدیک است.

## سوال ۴) جریان داده

نمودار خط :



❖ به نظر شما یک شرط تقریبی روی فرکانس واقعی عناصر که بتواند relative error را همواره

کمتر از یک نگه دارد، چیست؟

با توجه به رابطه‌ی relative error داریم:

$$Er[i] = (F'[i] - F[i]) / F[i]$$

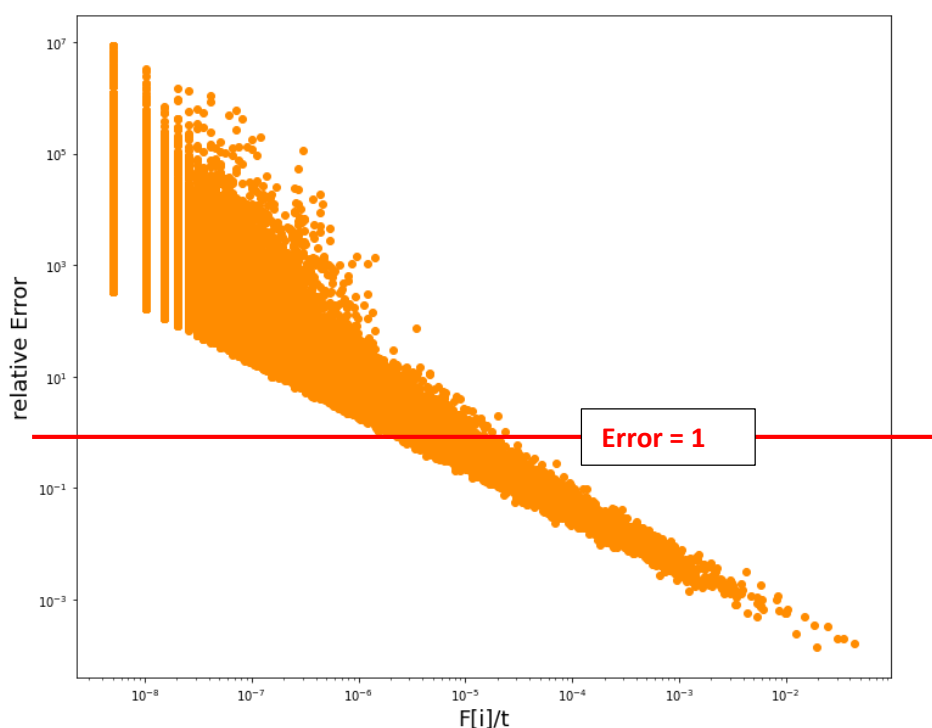
$$Er[i] < 1 \Rightarrow (F'[i] - F[i]) / F[i] < 1$$

$$F'[i] - F[i] < F[i] \Rightarrow F'[i] < 2 * F[i]$$

$$F[i] > F'[i]/2$$

بنابراین از روی رابطه به این شرط می‌رسیم که اگر فرکانس واقعی هر عنصر کوچکتر از نصف فرکانس تخمینی باشد خطا همواره کمتر از ۱ خواهد بود.

اگر بخواهیم از روی نمودار یک شرط پیدا کنیم خط  $error=1$  را رسم میکنیم :



با توجه به خط رسم شده، اگر مقدار فرکانس واقعی هر عنصر حدوداً بیشتر از  $10^{-5}$  باشد مقدار خطا کمتر از ۱ خواهد شد.