



BIG DATA

Homework #2



MALIHE GHAZALIAN

830402056
Spring 2024

سوال 1)

فیلتر بلوم با تعداد 1000 کد ملی معتبر که از دیتاست به صورت رندم انتخاب شد و 1000 کد ملی نامعتبر که به غیر از این 1000 تای اول بود از دیتاست سمپل شدند سپس با طول بیت 10000 و تعداد 10 تابع هش، الگوریتم پیاده سازی شد و تعداد کد ملی های معتبر و نامعتبر به صورت زیر تخمین زده شد:

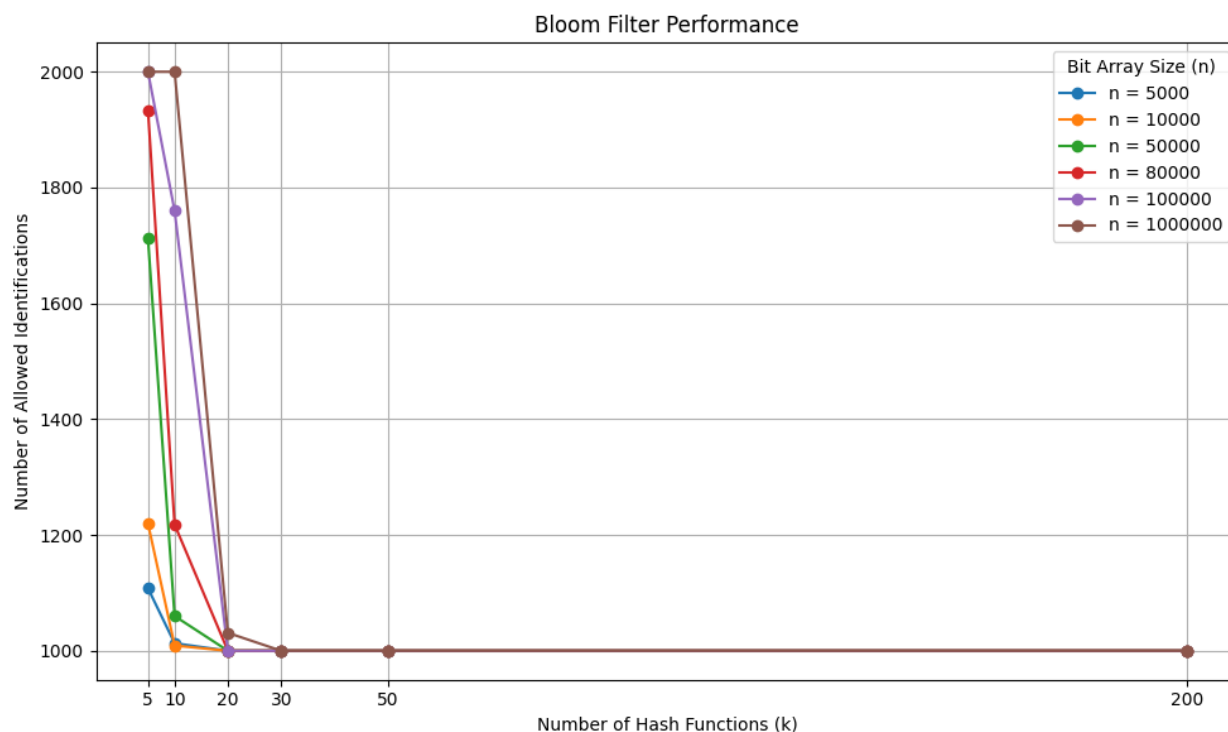
| Count of Allowed IDs: 1009

| Count of Not Allowed IDs: 991

در مرحله بعد با تعداد تابع هش متفاوت و طول های متفاوت بردار بیتی به شرح زیر الگوریتم اجرا شد :

```
num_hash_funcs = [5, 10, 20, 30, 50, 200]
bit_arr_sizes = [5000, 10000, 50000, 80000, 100000, 1000000]
```

نتایج تغییرات تعداد کدهای معتبر تشخیص داده شده به تعداد تابع هش در طول های متفاوت بردار بیتی به صورت نمودار زیر است:



بردار های بیتی در تعداد توابع کم تعداد false positive را نشان می دهد که بیان گر این است که استریم های جدید به خانه هایی هش می شوند که از قبل 1 بودند بدون آنکه عضو آن مجموعه

باشند ولی با افزایش تعداد هش و بیت های بردار بیتی این خطا کاهش می یابد و در تعداد 10 تا 20 هش و بردار بیتی با اندازه 5000 تا 10000 بهینه ترین جواب را می دهد.

اگر بخواهیم که احتمال fp برابر 0.01 باشد در فورمول احتمال یک بودن هر خانه که احتمال fp هم هست احتمال را قرار داده و از فرمول طول بیت را بدست می آوریم که فورمول به صورت زیر در می آید و با قرار دادن احتمال و تعداد المان ها که 1000 تا است:

$$M = -n \ln(P) / (\ln 2)^2$$

جوابش می شود : 9585.74 حدود 10000 هزار طول بردار بیتی باشد

سپس از فورمول k را استخراج می کنیم که به صورت : $k = n/m \ln(2)$ در می آید و با جایگذاری m می شود: 6.91 یعنی حدود 7 تابع هش نیاز است تا احتمال fp بشود 0.01

سپس با مقادیر فوق یکبار دیگر الگوریتم را اجرا کرده که به جواب زیر می رسیم که همان حدود 0.01 fp را می دهد:

Count of Allowed IDs: 1008 |

Count of Not Allowed IDs: 992 |

(سوال 2)

برای اینکه 35 تخمین داشته باشیم 35 تابع هش می خواهیم و با تخمین های مختلف بر اساس نتایج که به ازای هر استریم ورودی ثبت می شود روش ترکیبی نتایج به واقعیت نزدیک تری می دهد مثلا برای استریم اول روش ترکیبی 1 می شود ولی برای دو روش میانه و میانگین اعداد متفاوت است مخصوصا برای میانگین که کمی پرت می شود.

در قسمت بعد با طول های مختلف و با تعداد توابع مختلف الگوریتم اجرا شد و میانگین مقادیر پرتی را در استریم های بالاتر نشان می دهد همچنین با افزایش مقادیر میانگین بیشترین تغییرات را دارد.

سوال 3)

(الف)

به این علت که بعضی از آیت‌ها ممکن است در خیلی از بسکت‌ها ظاهر شوند و دلیل اصلی ظاهر شدن آیت‌ها دوم آیت‌ها اول نباشد به طوری که آیت‌ها دوم مستقل از آیت‌ها اول در بسیاری از بسکت‌ها ظاهر می‌شود بنابراین باید یک طوری این وابستگی را اثر دهیم یعنی تعداد بسکت‌هایی که آیت‌ها دوم در آن ظاهر شده اگر آیت‌های گروه ۱ در آن وجود دارند تقریباً برابر باشد بنابراین باید اثر ظاهر شدن کلی آیت‌ها که نتیجه می‌شود نیز ظاهر شود در دو متریک دیگر به دلیل اینکه در lift داریم مقدار کانفیدنس را به سائورت آیت‌ها دوم که به طور مستقل ظاهر می‌شود تقسیم میکنیم پس اگر آیت‌ها دوم خیلی ظاهر شده باشد به طور مستقل مقدار لیفت کم می‌شود در conv نیز به طور معکوس عمل می‌کنیم یعنی اگر تعداد باری که دومی ظاهر می‌شود کم باشد صورت زیاد می‌شود و مخرج هم که میتوانیم ثابت فرض کنیم زیرا میخواهیم اثر وجود آیت‌ها دوم به طور مستقل را بررسی کنیم پس با ثابت بودن کانفیدنس اگر آیت‌ها دوم کمتر به طور مستقل ظاهر شود مخرج بزرگ شده و مقدار conv بیشتر می‌شود.

(ب)

$$\textcircled{1} \text{ Confidence } \rightarrow (A \rightarrow B) = \frac{\text{Pr}(B|A)}{\text{support}(A)} = \frac{\text{support}(A, B)}{\text{support}(A)}$$

$$\textcircled{2} \text{ Lift } \rightarrow (A \rightarrow B) = \frac{\text{Conf}(A \rightarrow B)}{\text{support}(B)} = \frac{\text{support}(A, B) \times N}{\text{support}(A) \times \text{support}(B)}$$

$$\textcircled{3} \text{ Conviction } \rightarrow (A \rightarrow B) = \frac{1 - \text{support}(B)}{1 - \text{Conf}(A \rightarrow B)} = \frac{1 - \frac{\text{support}(B)}{N}}{1 - \frac{\text{support}(A, B)}{\text{support}(A)}}$$

$$\textcircled{1} \text{ Conf}(A \rightarrow B) = \frac{\text{support}(A, B)}{\text{support}(A)}$$

$$\text{Conf}(B \rightarrow A) = \frac{\text{support}(A, B)}{\text{support}(B)}$$

میزان درست بودن یک قاعده
 که از طریق نسبت درست به کل
 ظاهر شده است به این
 دو مقدار تقسیم شود
 در نتیجه یک معیار به دست می آید

$$\textcircled{2} \text{ Lift}(A \rightarrow B) = \frac{\text{support}(A, B) \times N}{\text{support}(A) \times \text{support}(B)}$$

$$\text{Lift}(B \rightarrow A) = \frac{\text{support}(A, B) \times N}{\text{support}(B) \times \text{support}(A)}$$

که هر ۲ مقدار برابر است پس
 این یک معیار در طرفه است

$$\textcircled{3} \text{ Conv}(A \rightarrow B) = \frac{1 - \frac{\text{support}(B)}{N}}{1 - \frac{\text{support}(A, B)}{\text{support}(A)}} = \frac{\frac{N - \text{support}(B)}{N}}{\frac{\text{support}(A) - \text{support}(A, B)}{\text{support}(A)}} = \frac{N - \text{support}(B)}{\text{support}(A) - \text{support}(A, B)}$$

$$\text{Conv}(B \rightarrow A) = \frac{1 - \frac{\text{support}(A)}{N}}{1 - \frac{\text{support}(A, B)}{\text{support}(B)}} = \frac{\frac{N - \text{support}(A)}{N}}{\frac{\text{support}(B) - \text{support}(A, B)}{\text{support}(B)}} = \frac{N - \text{support}(A)}{\text{support}(B) - \text{support}(A, B)}$$

$$\Rightarrow \text{Conv}(A \rightarrow B) = \frac{(N - \text{supp}(B)) \text{supp}(A)}{N(\text{supp}(A) - \text{supp}(A, B))}$$

$$\text{Conv}(B \rightarrow A) = \frac{(N - \text{supp}(A)) \text{supp}(B)}{N(\text{supp}(B) - \text{supp}(A, B))}$$

از نظر ریاضی به خاطر
تفاوت در صورت و مخرج
هر یک از این تعامیر ممکن

است با هم تفاوت می‌کند
نمیدان گفت که این یک معیار متفاوت است.

(ج)

Confidence صورت دلالت کامل، مقدارش برابر 1 است، بنابراین Confidence به حداکثر مقدار خود می‌رسد.

Lift در صورت دلالت کامل مقدار صورت که کانفیدنس است به یک میرسد ولی چون وابسته به ساپورت آیتم دوم است بنابراین به حداکثر مقدار نمی‌رسد و مطلوب نیست.

Conviction در صورت دلالت کامل چون مخرج برابر 1-1 می‌شود مقدارش برابر بینهایت است ولی اگر بینهایت را مقدار حداکثر آن در نظر بگیریم پس به حداکثر خود می‌رسد و مطلوب است

در قسمت بعدی سوال نیز الگوریتم apriori طبق خواسته های سوال اجرا شده و قوانین انجمنی برای جفت ها و سه تایی ها محاسبه شده و 5 تای برتر شناسایی شدند:

پنج قانون برتر دوتایی :

((('DAI93865', 'FRO40251'), 1.0), (('GRO85051', 'FRO40251'), 0.99918), (('GRO38636', 'FRO40251'), 0.99065), (('ELE12951', 'FRO40251'), 0.99057), (('DAI88079', 'FRO40251'), 0.98673))

پنج قانون برتر سه تایی:

((('DAI83733', 'ELE92920'), 'DAI62779'), 0.9279279279279279), ((('ELE17451',]
'ELE92920'), 'DAI62779'), 0.8984375), ((('ELE92920', 'DAI85309'), 'DAI62779'),
0.9502487562189055), ((('FRO85978', 'ELE59028'), 'DAI62779'), 0.874251497005988),
[(((('GRO81087', 'ELE92920'), 'DAI62779'), 0.9571428571428572)

سوال 4)

با همان دیتا الگوریتم های اصلاح شده ی a-priori اجرا شد که نتایج نشان می دهد که جوابهایی که الگوریتم ها به دست می دهند با هم متفاوت است مثلا در الگوریتم pcy با تغییر مقادیر سائز هش تعداد زوج ها تغییری میکند.