



الگوریتم
تعداد متمایز

محسن هوشمند
دانشکده علم رایانه و تکنولوژی اطلاعات
دانشگاه تحصیلات تکمیلی علوم پایه زنجان

تخمین تعداد منحصر به فرد

مسئله تخمین تعداد منحصر به فرد

- در پی یافتن تعداد مقادیر متمایز در مجموعه داده‌ای با وجود اعضای تکراری
- کاربرد در حوزه‌های گوناگونی از تولید نرم‌افزار اعم پایگاه داده‌ها و ترافیک شبکه

روش سنتی و دم دستی تعیین تعداد دقیق مجموعه

- ایجاد فهرستی از همه مقادیر
- اجتناب از شمارش عضو تکراری با مرتب‌سازی و جستجو
- احصاء تعداد مقادیر متمایز
- شمارش دقیق مقادیر منحصر به فرد با پیچیدگی زمانی $O(N \log N)$
- N تعداد کل مقادیر از جمله تکراری‌ها،
- نیاز به حافظه خطی کمکی

تخمین تعداد منحصر به فرد

گسترش خدمات اینترنتی،

- با میلیاردها کلیک، جستجو، و سفارش
- روزانه با تعداد بسیار کمتری کاربر
- اهمیت بازطراحی شمارش تعداد متمایز

تعیین تعداد منحصر به فرد

- با پیمایش یکباره داده و در فضای اشغالی کمتر
- از حوزه‌های مهم تحقیقاتی

تخمین تعداد منحصر به فرد

مثال‌های تخمین کاردینالیته

- تعیین تعداد بازدیدکننده یکتا از محصولی خاص
- تعداد کاربران متمایز استفاده کننده از قابلیتی مشخص در برنامه‌های وبی
- تشخیص تغییرات ناگهانی در تعداد آدرس‌های IP مبدأ—مقصد یکتا عبوری از روتر
- نشان‌دهنده امکان حمله انکار خدمت

تخمین تعداد منحصر به فرد

نحوه تکثیر مکرر اطلاعات در وب،

- اندازه‌گیری کاردینالیته کمک کننده در تشخیص مواجهه با تعداد محتوای متمایز
- مثال تعداد مقاله‌های خبری یکتا یا تعداد نسخه‌های متفاوت از محتوای وبسایت

داده‌مجموعه‌های عظیم امروز

- علاقه فزاینده‌ای به طراحی الگوریتم‌هایی جهت تقریب کاردینالیته مجموعه
- با دقتی مناسب و در فضایی بسیار کوچک‌تر از اندازه خود مجموعه

تخمین تعداد منحصر به فرد-مثال

از شاخص‌های کلیدی ارزشمند برای هر وبسایت

- تعداد بازدیدکنندگان منحصر به فردی بازدیدکننده در بازه زمانی مشخص
- جهت سادگی فرض استفاده هر بازدیدکننده از یک نشانی IP در تمامی بازدیدها
- با شمارش تعداد IP-های منحصر به فرد که با پروتکل اینترنت نسخه IPv6 با رشته‌های ۱۲۸ بیتی بدست آوردن تعداد بازدیدکنندگان
- آسان؟
- آیا می‌توان با تکیه صرف بر روش‌های کلاسیک برای شمارش دقیق تعداد را بدست آورد؟
- وابستگی پاسخ به میزان استفاده از تارمانه

تخمین تعداد منحصر به فرد-مثال

آمار ترافیک اسفند ۱۳۹۵ سه وبسایت خرده فروشی پر استفاده در ایالات متحده

▪ walmart.com و ebay.com، amazon.com

▪ مطابق [SimilarWeb](#)

▪ میانگین تعداد بازدید از این وبسایت‌ها حدود ۱,۴۴ میلیارد

▪ میانگین تعداد صفحات مشاهده شده در هر بازدید ۸,۲۴ بود. سخن کوتاه،

▪ آمار اسفند ۱۳۹۵ مراجعه به حدود ۱۲ میلیارد آدرس IP ۱۲۸۱ بی‌تی

▪ نگهداری همه آنها اندازه‌ای در حدود ۱۹۲ گیگابایت

▪ فرض به ازای هر ۱۰ بازدیدکننده، یک نفر منحصر به فرد

▪ انتظار مجموعه‌ای از تعداد منحصر به فرد حدود ۱۴۴ میلیون

▪ حافظه مورد نیاز برای ذخیره فهرست مقادیر منحصر به فرد ۲۳ گیگابایت

تخمین تعداد منحصر به فرد-مثال

مطالعه همبستگی‌ها در توالی‌های DNA

- از پژوهش‌های مهم در تحقیقات ژنوم انسانی

مولکول‌های DNA شامل دو رشته جفت شده

- هر کدام از چهار واحد شیمیایی پایه DNA، با علامت A (آدنین)، G (گوانین)، C (سیتوزین) و T (تیمین)

- ژنوم انسان حدود ۳ میلیارد جفت پایه

توالی‌یابی به معنای تعیین ترتیب دقیق جفت‌های پایه در قطعاتی از DNA

توالی DNA رشته‌ای از نمادهای A، G، C، T

- به اندازه دلخواه طولانی

- نمونه‌ای از مجموعه داده‌ای بالقوه بی‌نهایت

مسئله اندازه‌گیری همبستگی

- در پی تعیین تعداد زیررشته‌های متمایز با اندازه ثابت در قطعه‌ای DNA

با فرض

- هر توالی با زیررشته‌های متمایز کم، همبستگی بیشتری نسبت به توالی دیگری با اندازه مشابه اما با زیررشته‌های متمایز بیشتر

چنین آزمایش‌هایی نیاز به اجرای چندباره روی فایل‌های بزرگ متعدد

نیاز به حافظه محدود یا حتی ثابت و زمان اجرای کم برای تسریع تحقیقات

- غیرقابل اجرا با الگوریتم‌های شمارش دقیق

تخمین تعداد منحصر به فرد

دستاوردهای احتمالی تخمین دقیق تعداد منحصر به فرد نیاز زمان پردازش و حافظه زیاد

گرایش برنامه‌های کلان‌داده به سمت رویکردهای عملی‌تر
▪ غالباً بر مبنای الگوریتم‌های احتمالاتی و عرضه پاسخی تقریبی

نیاز به داشتن درکی از اندازه مجموعه داده و تعداد مقادیر متمایز احتمالی در تصمیم‌گیری

▪ دنباله بی‌نهایت (بالقوه بی‌نهایت) رشته‌های تک حرفی a، d، s، ... بر اساس حروف الفبای لاتین
▪ راحتی امکان تخمین تعداد منحصر به فرد

▪ کران بالای آن با تعداد حروف، که در الفبای لاتین فعلی در انگلیسی ۲۶

▪ عدم نیاز به استفاده از هیچ رویکرد احتمالاتی

▪ راه‌حل مبتنی بر لغت‌نامه ساده برای محاسبه دقیق تعداد منحصر به فرد

تخمین تعداد منحصر به فرد

برای پرداختن به مسئله تعداد منحصر به فرد،

تاثیر بسیاری از روش‌های احتمالاتی مطرح از الگوریتم فیلتر بلوم

عمل بر اساس مقادیر درهم

مشاهده الگوهای رایج در توزیع آنها

«حدس‌های» منطقی در مورد تعداد مقادیر منحصر به فرد بدون نیاز به ذخیره همه آنها

شمارش اقلام متمایز در پایگاه‌های داده

برای سنجش کاردینالیته (اندازهٔ مجموعه) از پایگاه‌های داده و نحوهٔ استفادهٔ SQL از کلیدواژهٔ DISTINCT

اعمال کلیدواژه بر روی یک ستون در جدول

- دستور SELECT DISTINCT بازگرداندن تمام اقلام متمایز آن ستون
- دستور SELECT COUNT DISTINCT برگرداندن تعداد اقلام متمایز موجود در ستون مورد نظر
- رواج پرس‌وجوهای حاوی COUNT DISTINCT
 - به‌ویژه در حوزهٔ تجارت الکترونیک
 - آمار استفاده از وبسایت
- ذخیرهٔ داده‌های بازدید کاربران اغلب در جدولی به نام DAILY_VISITS
 - معمولا بسیار بزرگ
- دارای ویژگی‌هایی نظیر session_id, timestamp, product_id, user_ip_address, visit_duration و غیره
- اجرای عملیات SELECT زیر:

```
SELECT COUNT (DISTINCT user_ip_address) WHERE product_id = #####  
FROM DAILY_VISITS
```

شمارش ارقام متمایز در پایگاه‌های داده

دریافت تعداد آدرس‌های IP متمایز (یعنی کاربران) دسترسی یافته به محصول با شناسه ##### در یک روز معین

- در وبسایت شلوغ، رشد جدول بازدیدهای روزانه تا چند میلیارد ردیف
- زمان زیاد اجرا و تاخیر این پرس‌وجوی

شمارش ارقام متمایز در پایگاه‌های داده

تأخیر عمدتاً ناشی از عملیات مرتب‌سازی

- **COUNT DISTINCT** کلاسیک در بیشتر پایگاه‌های داده (مانند Azure SQL/SQL Server)
- مگر آنکه ستون از قبل مرتب
- پس از مرتب‌سازی ستون، قرارگیری تمام رکوردهای تکراری در کنار یکدیگر
- کفایت پیمایش ترتیبی جهت شناسایی و شمارش ارقام متمایز
- هزینه عملیات مرتب‌سازی در جدولی با n ردیف، $O(n \log n)$
- مقیاس‌ناپذیری نه تنها برای چند میلیارد ردیف بلکه حتی برای چند میلیون ردیف
- پرس‌وجوهای ساده نیز دارای **COUNT DISTINCT**‌ها و **GROUP BY**‌های متعدد روی ستون‌های مختلف
- عدم تاثیر مرتب‌سازی یک ستون بر کاهش پیچیدگی مرتب‌سازی ستون دیگر
- جهت تسریع کار استفاده از جدول درهم
- نیاز جدول درهم به فضایی خطی متناسب تعداد مقادیر متمایز k
- امکان افزایش k تا n
- مقرون به صرفه نبودن استفاده از درهم

شمارش ارقام متمایز در پایگاه‌های داده

نیاز صرف به دانستن تعداد ارقام متمایز و عدم نیاز به فهرست کردن خود ارقام متمایز

- پابرجایی پیچیدگی
- مسئله متمایز بودن مقادیر (element-distinctness)
- با داشتن آرایه‌ای از n عنصر،
- تعیین متمایز بودن تمام مقادیر آن
- دارای کران پایین $\Omega(n \log^2 n)$

شمارش ارقام متمایز در پایگاه‌های داده

برای رفع مسائل مقیاس‌پذیری

نسخه‌های جدیدتر سیستم‌های مدیریت پایگاه‌داده و انبارهای داده

- تخمین کاردینالیته روی آورده‌اند: *SQL Server 2019*

- دارای عملیات `APPROX_COUNT_DISTINCT`

- استفاده از فضای بسیار کمی و سرعت بالا

- **Google BigQuery**

- استفاده از این رویکرد تقریبی و احتمالاتی در `COUNT DISTINCT` به عنوان پیش‌فرض

- استفاده از `EXACT_COUNT_DISTINCT` برای موقعیت‌هایی که مطلقاً به پاسخ دقیق نیاز است

- استفاده از ابرلگ در این تخمین‌ها

- ابداع فلاژوله و همکاران

- صرفه‌جویی بالا در فضا (در حد کیلوبایت) با وجود کار روی مجموعه‌داده‌هایی در حد تریلیون

- میزان خطا نسبتاً پایین

- در مرتبه‌ای $O(\frac{1}{\sqrt{m}})$ که m نشان‌دهنده تعداد شمارنده‌ها یا بیت-نگاشت‌های چند بیتی

- انتخاب رایج برای m مقدار 2^{14}

شمارش ارقام متمایز در پایگاه‌های داده

- مشاهده نمونه‌های متعددی از صرفه‌جویی در فضا به بهای کاهش دقت
- معنابخشی دیگر ابرلگ‌لگ به کارایی فضا
 - باقیماندن تقریباً همیشه در محدوده چند کیلوبایت
 - به دست دادن کاردینالیت‌ها واقعی با میانگین نسبت خطای کوچک (مثلاً ± 2 درصد)

شمارش خطی

الگوریتم شمارش خطی

- اولین رویکرد احتمالاتی برای مسئله تعداد منحصر به فرد
- مورتون آستراهان، ماریو شولنیک و کیو-یانگ وانگ در سال ۱۳۶۶
- پیشنهاد الگوریتم عملی را کیو-یانگ وانگ، براد واندر-زاندن، و هوارد تیلور در سال ۱۳۶۹

طرح اصلی پیشنهادی آنها

- موجب برتری نسبت به روش‌های دقیق
 - استفاده درهم مقادیر با تابعی درهم‌ساز
 - فراهم شدن امکان حذف تکراری‌ها بدون نیاز به مرتب‌سازی مقادیر
 - دارای هزینه نیست: کاهش دقت پاسخ
 - احتمال خطا به دلیل تصادم‌های احتمالی درهم
 - عدم امکان تشخیص تکراری‌ها خاصه «تکراری‌های تصادفی»
- استفاده از چنین جدول درهم به همراه روش پیمایش مناسب
- منجر به بهبود عملکرد فضا و زمان نسبت به روش کلاسیک

شمارش خطی

برای مجموعه داده‌های با تعداد منحصر به فرد عظیم

- امکان بزرگ بودن چنین جدول‌های درهم
- نیاز به حافظه‌ای با رشد خطی با تعداد مقادیر متمایز در مجموعه

برای سیستم‌هایی با حافظه محدود

- نیاز به ذخیره‌سازی دیسک یا توزیع شده
- به دلیل دسترسی کند به دیسک یا شبکه کاهش مزایای جدول‌های درهم

شمارش خطی

مشابه ایده فیلتر بلوم برای دور زدن چنین مشکلی

الگوریتم شمارش خطی

- ذخیره نکردن خود مقادیر درهم
- ذخیره بیت‌های متناظر آنها
- جانشینی جدول درهم با آرایهٔ بیتی (تحت «شمارندهٔ خطی») به طول m
- فرض بر تناسب m با تعداد مقادیر متمایز مورد انتظار n
- اما نیاز به یک بیت به ازای هر مقدار برای اکثر موارد

شمارش خطی

در ابتدا، همه بیت‌ها در «شمارنده خطی» برابر با صفر

درج مقدار جدید x

▪ محاسبه مقدار درهم‌ساز $h(x)$ متناظر

▪ برابر یک قرار دادن مقدار بیت در نمایه متناظر در شمارنده

الگوریتم

الگوریتم ۱- درج در شمارنده خطی

ورودی: مقدار $x \in D$

ورودی: شمارنده خطی با تابع درهم‌ساز h

$j \leftarrow h(x)$

اگر $LINEARCOUNTER[j] = 0$ آنگاه

$LINEARCOUNTER[j] \leftarrow 1$

شمارش خطی

به دلیل استفاده از تک تابع درهم‌ساز h

▪ امکان تصادم

▪ دو مقدار ورودی متفاوت منجر به مقداردهی یک بیت در آرایه

▪ در نتیجه تعداد تصادم‌های فراوان

▪ عدم امکان به دست آوردن مستقیم تعداد دقیق (یا حتی نزدیک به دقیق) مقادیر متمایز از چنین طرح فشرده‌ای

شمارش خطی

- ایده الگوریتم منجر به توزیع مقادیر در سطرها
- بیت‌های نمایه شده مقادیر درهم
 - آرایهٔ بیتی نمایشگر مقداردهی شدن مدخل‌ها
 - مشاهده تعداد یک‌ها در آرایه تخمین تعداد منحصربه‌فرد

شمارش خطی

در مرحله اول الگوریتم شمارش خطی،
▪ ایجاد داده ساختار

با داشتن چنین طرح فشرده‌ای

▪ تخمین تعداد منحصربه‌فرد با استفاده از کسر مشاهده شده بیت‌های خالی V به صورت زیر

$$n \approx -m \cdot \ln V$$

شمارش خطی

تاثیر تصادم‌ها بر تخمین تعداد منحصربه‌فرد در الگوریتم شمارش خطی با هر تصادم کاهش تعداد بیت‌هایی که باید تنظیم شوند کسر مشاهده شده بیت‌های تنظیم نشده بزرگتر از مقدار واقعی نداشتن تصادم درهم

- تعداد نهایی بیت‌های تنظیم شده برابر تعداد منحصربه‌فرد مورد نظر
- اما اجتناب‌ناپذیری تصادم‌ها
- معادله به مثابه تخمینی بیش از حد از تعداد دقیق منحصربه‌فرد

مقدار صحیح بودن تعداد منحصربه‌فرد

- گرد کردن نتیجه آن به نزدیک‌ترین عدد صحیح کوچکتر

شمارش خطی

الگوریتم ۲- تخمین تعداد منحصر به فرد با شمارش خطی

ورودی: مجموعه داده D

خروجی: تخمین تعداد منحصر به فرد

$LINEARCOUNTER[i] \leftarrow 0, i = 0 \dots m - 1$

برای هر $x \in D$ انجام بده

$Add(x, LINEARCOUNTER)$

$Z \leftarrow count_{i=0 \dots m-1} (LINEARCOUNTER[i] = 0)$

برگرداندن $\left\lceil -m \times \ln \frac{Z}{m} \right\rceil$

شمارش خطی-مثال

مجموعه داده‌ای شامل بیست نام از پایتخت‌ها

▪ برلین، برلین، پاریس، برلین، لیسبون، کیف، پاریس، لندن، رم، آتن، مادرید، وین، رم، لیسبون، برلین، پاریس، لندن، کیف، واشنگتن

▪ تعداد منحصر به فرد کوچک

▪ تعداد دقیق منحصر به فرد ۱۰ است

برای داشتن خطای استاندارد حدود ده درصد

▪ نیاز به انتخاب طول داده‌ساختار حداقل به اندازه تعداد مقادیر منحصر به فرد مورد انتظار است

$$m = 2^4$$

▪ تابع درهم‌ساز h با مقادیر در $\{0, 1, \dots, 2^4 - 1\}$ از تابعی مبتنی بر مورمور ۳ نوع ۳۲ بیتی

$$h(x) = \text{MurmurHash3}(x) \% m$$

شمارش خطی-مثال

مقادیر درهم

شهر	h (شهر)	شهر	h (شهر)
آتن	۱۲	لندن	۱۴
برلین	۷	لیسبون	۱۵
پاریس	۸	مادرید	۱۴
رم	۱	واشنگتن	۱۱
کیف	۱۳	وین	۶

شمارش خطی-مثال

شهرهای لندن و مادرید دارای مقداری یکسان
انتظار چنین تصادم‌هایی مورد انتظار و کاملاً طبیعی

۰	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱	۱۲	۱۳	۱۴	۱۵
۰	۱	۰	۰	۰	۰	۱	۱	۱	۰	۰	۱	۱	۱	۱	۱

شمارش خطی-مثال

با توجه به الگوریتم شمارش خطی
▪ محاسبه کسر V بیت‌های خالی

$$V = \frac{7}{16}$$

تعداد منحصر به فرد تخمین زده شده برابر

$$n \approx -16 \times \ln \frac{7}{16} \approx 13$$

شمارش خطی-ویژگی‌ها

محاسبه V

احتمال صفر بودن بیتی پس از درج n عضو جدید برابر است با

$$\text{pr}(\text{bit} = 0) = \left(1 - \frac{1}{m}\right)^n \approx e^{-n/m}$$

در نتیجه، میانگین نسبت یک‌ها برابر است با

$$\frac{E[N]}{m} = 1 - e^{-n/m}$$

و N تعداد بیت‌های متناظر یک در فیلتر

جهت تقریب n از N می‌توان از $N/m = 1 - e^{-n/m}$ مقدار زیر را بدست آورد:

$$e^{-n/m} = 1 - \frac{N}{m}$$

شمارش خطی-ویژگی ها

$$\begin{aligned}e^{-n/m} &= 1 - \frac{N}{m} \\ -\frac{n}{m} &= \ln \left(1 - \frac{N}{m} \right) \\ n &= -m \ln \left(1 - \frac{N}{m} \right), V = 1 - \frac{N}{M} \\ &\Rightarrow n = -m \ln V\end{aligned}$$

شمارش خطی-ویژگی ها

در صورتی که تابع درهم‌ساز h دارای زمان محاسبه‌ای ثابت
▪ صادق برای غالب توابع درهم‌ساز

زمان پردازش هر مقدار برابر مرتبه $O(N)$
▪ N تعداد کل مقادیر از جمله تکراری‌ها.
▪ الگوریتم دارای پیچیدگی زمانی $O(N)$

شمارش خطی-ویژگی ها

مانند هر الگوریتم احتمالاتی دیگر

- دارای تعدادی پارامتر قابل تنظیم جهت تاثیر بر عملکرد آن

دقت مورد انتظار تخمین

- به اندازه آرایهٔ m بیتی

- نسبت آن به تعداد مقادیر متمایز $\alpha = \frac{m}{n}$

- $\alpha \geq 1$

- $m > n$ که در عمل جالب نیست

احتمالی غیرصفر اشباع pr

- پر شدن آرایهٔ بیتی یا احتمال اشباع

- مخدوش سازی الگوریتم و بی نهایت کردن مقدار معادله

- وابستگی احتمال اشباع pr به ضریب بار و در نتیجه به اندازه m

- انتخاب به اندازهٔ کافی بزرگ

- جهت ناچیز کردن احتمال اشباع

شمارش خطی

خطای استاندارد δ

معیاری از درستی تخمین حاصل از شمارش خطی

تبادل بین آن و اندازه آرایه بیتی m

کاهش خطای استاندارد منجر به تخمین‌های دقیق‌تر، اما افزایش حافظه موردنیاز

جدول ۳. تبادل بین دقت و اندازه آرایه بیتی

n	m	
	$\delta = 1\%$	$\delta = 10\%$
۱۰۰۰	۵۳۲۹	۲۶۸
۱۰۰۰۰	۷۹۶۰	۱۷۰۹
۱۰۰۰۰۰	۲۶۷۲۹	۱۲۷۴۴
۱۰۰۰۰۰۰	۱۵۴۱۷۱	۱۰۰۸۸۰
۱۰۰۰۰۰۰۰	۱۰۹۶۵۸۲	۸۳۱۸۰۹
۱۰۰۰۰۰۰۰۰	۸۵۷۱۰۱۳	۷۰۶۱۷۶۰

شمارش خطی

وابستگی به انتخاب m بسیار پیچیده و نداشتن راه حل تحلیلی
نویسندگان الگوریتم مقادیر از پیش محاسبه شده‌ای برای احتمال قابل قبول تکمیل

پیشنهاد $pr = 0.7\%$ اشباع

- استفاده به عنوان مرجع
- به دلیل صفر نبودن احتمال تکمیل، نادر بودن پر شدن آرایه بیتی
- هنگام کار با مجموعه‌های داده کوچک،
- امکان نمایه‌سازی دوباره همه مقادیر با تابع درهم‌ساز متفاوت
- یا افزایش اندازه شمارنده خطی
- ناکارآمدی چنین راه‌حلی برای مجموعه‌های داده عظیم
- همراه با پیچیدگی زمانی نسبتاً بالا

شمارش خطی

مناسب بودن شمارش خطی

- بزرگ نبودن تعداد منحصربه‌فرد مجموعه داده‌ای
- امکان استفاده از آن در بهبود الگوریتم‌های دیگر برای تعداد منحصربه‌فردهای عظیم
- تخمین تعداد منحصربه‌فرد تقریباً متناسب با مقدار دقیق
- دلیل اصطلاح «خطی»

شمارش عناصر منحصر به فرد در فیلتر بلوم

$$e^{-kn/m} = 1 - \frac{N}{m}$$
$$-\frac{kn}{m} = \ln\left(1 - \frac{N}{m}\right)$$
$$n = -\frac{m}{k} \ln\left(1 - \frac{N}{m}\right)$$

که دقیقا برابر مقداری است که الگوریتم وقتی $0 < N < m$ برمی گرداند. حالت های خاص آن عبارت است از الف- $N < k$ که غیرممکن است مگر $n = 0$. ب- $N = k$ که نشان می دهد احتمالا یک عضو افزوده شدن. ج- $N = m$ یا تمامی بیت ها برابر ۱ است. در نتیجه ظرفیت بیشینه یا m/k عضو دارد.

معادله $-\frac{m}{k} \ln(1 - N/m)$ تقریبی مناسب از n در مواقعی است که فیلتر اشباع نشده است.

شمارش احتمالاتی

فیلیپ فلاژوله و جی. نایجل مارتین

▪ سال ۱۳۶۴

معرفی الگوریتم شمارشی بر اساس مشاهده الگوهای متداول در نمایش‌های درهم شده مقادیر نمایه شده

اعمال تابع درهم‌ساز h بر مقدار ورودی

تبدیل آن به عددی از اعداد صحیح با توزیع نزدیک-یکنواخت در بازه $\{0, 1, \dots, 2^\ell - 1\}$

▪ یا به طور معادل، تبدیل به رشته‌ای از مجموعه رشته‌های دودویی به طول ℓ

$$h(x) = i = \sum_{j=0}^{\ell-1} 2^j \times i_j = (i_{\ell-1} \dots i_1 i_0)_2, i_j \in \{0, 1\}$$

شمارش احتمالاتی

فلاژوله و مارتین تشخیص دادند

- هم احتمال بودن حضور الگوهای $10^j = \overbrace{100\dots0}^j$ در رشته‌های دودوئی وجود چنین مقداری در رشته دودوئی با احتمال $2^{-(j+1)}$
- در صورت ثبت شدن برای هر مقدار نمایه شده
- امکان ایفای نقش تخمین‌گر تعداد منحصر به فرد
- مرتبط کردن هر الگو با شاخص آن، به نام «رتبه»:

$$\text{رتبه}(i) = \begin{cases} \min j, \forall i > 0 \\ \ell, & i = 0 \end{cases}$$

- اساساً متناظر موقعیت سمت راست‌ترین ۱
- با کمترین اهمیت

شمارش احتمالاتی - مثال

عدد صحیح ۸ بیتی ۴۲ با:

$$\begin{aligned} & 42 \\ = & 0 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 \\ = & (00101010)_2 \end{aligned}$$

ظاهر شدن یک‌ها در موقعیت‌های ۱، ۳ و ۵. با توجه به تعریف، رتبه ۴۲ برابر با
 $\text{رتبه}(42) = \min(1,3,5) = 1$

شمارش احتمالاتی

- رخدادهای الگوی 10^j ،
- یا متناظرا $۱ + j = \text{رتبه}(\cdot)$ ،
- نمایه شده در نمایش‌های دودویی مقادیر درهم هر مقدار
- امکان ذخیره به صورت فشرده در داده‌ساختار آرایه‌ای شمارنده به طول l
- مشهور به «طرح فم»

در ابتدا، همه بیت‌ها در شمارنده برابر با صفر

هنگام درج مقدار جدید x در داده‌ساختار،

محاسبه مقدار درهم آن با تابع درهم‌ساز h

محاسبه رتبه متناظر

برابر یک کردن بیت متناظر در آرایه

شمارش احتمالاتی

الگوریتم ۳: درج مقدار در شمارنده ساده

ورودی: مقدار $x \in D$

ورودی: شمارنده ساده با تابع درهم‌ساز h

$j \leftarrow \text{رتبه}(h(x))$

اگر $COUNTER[j] == 0$ آنگاه

$COUNTER[j] \leftarrow 1$

شمارش احتمالاتی

یک در شمارنده در موقعیت $j + 1$
به معنای حداقل یک بار مشاهده الگوی 10^j در میان مقادیر درهم همه ورودی‌های نمایه شده

شمارش احتمالاتی

الگوریتم ۴- محاسبه رتبه X

ورودی: عدد صحیح X

خروجی: رتبه X یا کوچکترین j که بیت j-ام برابر ۱

$j \leftarrow 0$

$l = \text{len}(x)$

تازمانی که $j < l$ ادامه بده

اگر $(x \gg j \& 1) == 1$

برگرداندن j

$j \leftarrow j + 1$

برگرداندن 1 // عملاً رخ نمی‌دهد چون $h(x)$ غیرصفر

شمارش احتمالاتی - مثال

مجموعه داده بیست نام از پایتخت‌ها اعم از

برلین، برلین، پاریس، برلین، لیسبون، کیف، پاریس، لندن، رم، آتن، مادرید، وین، رم، رم، لیسبون، برلین، پاریس، لندن، کیف، واشنگتن

از مورمور^۳-^{۳۲} بیتی به عنوان تابع درهم‌ساز h

نگاشت به $\{0, 1, \dots, 2^{32} - 1\}$

استفاده از شمارنده‌ای به طول $\ell = 32$

با استفاده از مقادیر درهم

محاسبه رتبه‌های هر مقدار

شمارش احتمالاتی-مثال

رتبه	h (شهر)	شهر	رتبه	h (شهر)	شهر
۱	۳۴۵۰۹۲۷۴۲۲	لندن	۲	۴۱۶۱۴۹۷۸۲۰	آتن
۰	۶۲۹۵۵۵۲۴۷	لیسبون	۰	۳۶۸۰۷۹۳۹۹۱	برلین
۱	۲۹۷۰۱۵۴۱۴۲	مادرید	۳	۲۶۷۳۲۴۸۸۵۶	پاریس
۰	۴۰۳۹۷۴۷۹۷۹	واشنگتن	۰	۵۰۱۲۲۷۰۵	رم
۱	۳۲۷۱۰۷۰۸۰۶	وین	۰	۳۴۹۱۲۹۹۶۹۳	کیف

شمارش احتمالاتی

با فرض توزیع یکنواخت مقادیر

فرض تعداد مقادیر متمایز نمایه شده تاکنون برابر n

ظاهر شدن مقدار ۱ در مدخل اول در حدود $\frac{n}{2}$ دفعه

ظاهر شدن مقدار یک در مدخل دوم در حدود $\frac{n}{2^2}$ دفعه

ادامه به همین منوال

ظهور یک در موقعیت j برابر $\frac{n}{2^j}$ دفعه

شمارش احتمالاتی

اگر $j \gg \log_2 n$

- احتمال مشاهده مقدار ۱ در موقعیت j -ام میل به صفر
- در نتیجه مقدار $Counter[j]$ قریب به یقین برابر صفر

اگر $j \ll \log_2 n$

- مقدار $Counter[k]$ قریب به یقین برابر با یک

اگر مقدار $j \approx \log_2 n$

- برابر بدون احتمال مشاهده ۱ یا ۰ در آن موقعیت

شمارش احتمالاتی

سخن کوتاه

- پس از درج تمام مقادیر مجموعه داده در آرایه شمارنده
- سمت راست‌ترین موقعیتی با مقدار صفر ظاهر می‌شود

▪ به عنوان شاخص تخمین $\log_2 n$

- فلاژوله و مارتین در عمل، نیاز به اعمال ضرب تصحیح ϕ جهت بهبود تخمین

تخمین کاردینالیته (تعداد مقادیر متمایز)

$$n \approx \frac{1}{\phi} 2^R$$

▪ $\phi \approx 0.77351$

شمارش احتمالاتی

سخن کوتاه

- استفاده فلاژوله و مارتین از موقعیت کم‌ارزش‌ترین بیت صفر (چپ‌ترین صفر در آرایه) به عنوان تخمین کاردینالیته
- تعریف الگوریتم بر همین اساس

با توجه به مشاهده قبلی

- امکان استفاده از موقعیت پرارزش‌ترین بیت ۱ (راست‌ترین ۱) به همین منظور
- اما دارای توزیع یکنواخت‌تر و در نتیجه خطای استاندارد بزرگ‌تر

شمارش احتمالاتی

الگوریتم ۵- محاسبه موقعیت سمت چپ‌ترین صفر
ورودی: شمارنده ساده به طول l
خروجی: سمت چپ‌ترین موقعیت صفر
برای j از 0 تا $l - 1$ انجام بده
اگر $COUNTER[j] == 0$ برگرداندن j
برگرداندن l

شمارش احتمالاتی-مثال

۰	۱	۲	۳	۴	۵	۶	۷	۸	۹	۱۰	۱۱	۱۲	۱۳	۱۴	۱۵
۱	۱	۱	۱	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰
۱۶	۱۷	۱۸	۱۹	۲۰	۲۱	۲۲	۲۳	۲۴	۲۵	۲۶	۲۷	۲۸	۲۹	۳۰	۳۱
۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰	۰

با استفاده از الگوریتم 5، در شمارنده ظاهر شدن چپ‌ترین مقدار ۰ در موقعیت $R = 4$ تخمین تعداد منحصر به فرد

$$n \approx \frac{1}{0.77351} 2^4 \approx 20.68$$

- تعداد دقیق منحصر به فرد مجموعه ۱۰
- خطای زیاد تخمین
- دلیل عدد صحیح بودن مقادیر R
- برای رتبه‌های بسیار نزدیک
- امکان حصول نتایج دارای تفاوتی از مرتبه توان‌های دودوئی
- در مثال حاضر، $R=3$ برگرداندن تخمین نسبتاً مناسب ۱۰,۳۴

شمارش احتمالاتی

تخمین تعداد منحصربه‌فرد بر اساس شمارنده معمول
▪ بدست آوردن مقادیر میانگین بسیار نزدیکی، اما با وردایی نسبتاً بالا
▪ در مثال قبل، با خطای استاندارد δ دارای اختلاف از مرتبه دو برابر

ضعف رویکرد تک‌شمارنده
▪ فقدان تخمین‌های با اطمینان بالا برای تعداد منحصربه‌فرد
▪ انجام پیش‌بینی تنها بر اساس یک تخمین

بهبود منطقی
▪ استفاده از تعداد زیادی شمارنده
▪ منجر به افزایش تعداد تخمین‌ها

▪ پیش‌بینی نهایی n با میانگین‌گیری از پیش‌بینی‌های R_i حاصل شمارنده‌های مزبور $\{ \text{شمارنده}_i \}_{i=0}^{m-1}$

سخن کوتاه، معادله اصلاح شده الگوریتم شمارش احتمالاتی:

$$n \approx \frac{1}{\phi} 2 \frac{1}{m} \sum_{i=0}^{m-1} R_i$$

شمارش احتمالاتی

تعداد منحصر به فرد n

- مقدار تخمینی با همان کیفیت اما با وردایی بسیار کوچکتر
- اشکال عملی برای ساخت m شمارنده ساده مستقل،
 - نیاز به محاسبه مقادیر m تابع درهم ساز متفاوت
 - محاسبه هر تابع درهم ساز در $O(1)$
- محاسبه همه آنها دارای پیچیدگی زمانی $O(m)$ و هزینه‌های پردازنده نسبتاً بالا

شمارش احتمالاتی

راه حل بهینه‌سازی الگوریتم شمارش احتمالاتی

- اعمال «میانگین‌گیری تصادفی»
- جانشینی تمامی m تابع درهم‌ساز با صرفاً یک تابع
- تعریف m شمارنده
- تقسیم مقدار حاصل به دو بخش خارج‌قسمت و باقیمانده
- باقیمانده r برای انتخاب یکی از m شمارنده و خارج‌قسمت q برای محاسبه رتبه
- یافتن شاخص مناسب برای به‌روزرسانی در آن شمارنده

شمارش احتمالاتی

الگوریتم ۶: استفاده از میانگین‌گیری برای به‌روزرسانی شمارنده‌ها

ورودی: مقدار $x \in D$

ورودی: آرایه‌ای از m شمارنده ساده با تابع درهم‌ساز h

$$f_r \leftarrow h(x) \% m$$

$$f_q \leftarrow \left\lfloor \frac{h(x)}{m} \right\rfloor$$

$$j \leftarrow \text{rank}(f_q)$$

$$\text{COUNTER}_{f_r}[j] \leftarrow 1$$

شمارش احتمالاتی

اعمال الگوریتم میانگین گیری تصادفی به شمارش احتمالاتی،
▪ با فرض کفایت منصفانه بودن توزیع مقادیر مبتنی بر خارج قسمت

انتظار نمایه شدن $\frac{n}{m}$ از کل مقدار در هر شمارنده $\left\{ \text{شمارنده } i \right\}_{i=0}^{m-1}$

در نتیجه: معادله قبلی تخمین خوبی برای $\frac{n}{m}$

▪ نه مستقیماً n

▪ در نتیجه پیشنهاد فلاژوله و مارتین

$$n \approx \frac{m}{\phi} 2^{\bar{R}} = \frac{m}{\phi} 2^{\frac{1}{m} \sum_{i=0}^{m-1} R_i}$$

شمارش احتمالاتی

الگوریتم ۷: الگوریتم فلاژوله-مارتین (PCSA)

ورودی: مجموعه داده D

ورودی: آرایه‌ای از m شمارنده ساده با تابع درهم‌ساز h

خروجی: تخمین تعداد منحصربه‌فرد

برای $x \in D$ انجام بده

$$f_r \leftarrow h(x) \% m$$

$$f_q \leftarrow \left\lfloor \frac{h(x)}{m} \right\rfloor$$

$$j \leftarrow \text{rank}(f_q)$$

$$\text{COUNTER}_{f_r}[j] \leftarrow 1$$

برای i از 0 تا $m - 1$ انجام بده

$$R \leftarrow \text{LeftMostZero}(\text{COUNTER}_{f_r})$$

$$S \leftarrow S + R$$

$$S \leftarrow \frac{1}{m} S$$

برگرداندن $\frac{m}{\phi} 2^S$

شمارش احتمالاتی

الگوریتم شمارش احتمالاتی با میانگین گیری تصادفی (PCSA)

- همچنین الگوریتم فلاژوله-مارتین
- در مقایسه با نسخه‌ای که از m تابع درهم‌ساز
- کاهش پیچیدگی زمانی به حدود $O(1)$

شمارش احتمالاتی-مثال

تخمین تعداد منحصر به فرد با میانگین گیری تصادفی

اعمال میانگین گیری تصادفی با شبیه سازی $m = 3$ تابع درهم ساز

استفاده از باقیمانده r برای انتخاب یکی از سه شمارنده و از خارج قسمت q برای محاسبه رتبه

رتبه (q)	q	r	شهر (h)	شهر	رتبه (q)	Q	r	شهر (h)	شهر
۲	۱۱۵۰۳۰۹۱۴۰	۲	۳۴۵۰۹۲۷۴۲۲	لندن	۲	۱۳۸۷۱۶۵۹۴۰	۰	۴۱۶۱۴۹۷۸۲۰	آتن
۰	۲۰۹۸۵۱۷۴۹	۰	۶۲۹۵۵۵۲۴۷	لیسبون	۱	۱۲۲۶۹۳۱۳۳۰	۱	۳۶۸۰۷۹۳۹۹۱	برلین
۲	۹۹۰۰۵۱۳۸۰	۲	۲۹۷۰۱۵۴۱۴۲	مادرید	۳	۸۹۱۰۸۲۹۵۲	۰	۲۶۷۳۲۴۸۸۵۶	پاریس
۰	۱۳۴۶۵۸۲۶۵۹	۲	۴۰۳۹۷۴۷۹۷۹	واشنگتن	۴	۱۶۷۰۷۵۶۸	۱	۵۰۱۲۲۷۰۵	رم
۰	۱۰۹۰۳۵۶۹۳۵	۱	۳۲۷۱۰۷۰۸۰۶	وین	۲	۱۱۶۳۷۶۶۵۶۴	۱	۳۴۹۱۲۹۹۶۹۳	کیف

شمارش احتمالاتی-مثال

هر شمارنده مدیریت اطلاعات مربوط به حدود یک سوم شهرها

توزیع به اندازه کافی منصفانه

پس از نمایه‌سازی همه مقادیر و تنظیم بیت‌های مناسب در شمارنده‌های مربوطه،

به دست آوردن شمارنده‌ها

شمارش احتمالاتی-مثال

سمت چپ‌ترین موقعیت‌های صفر برای هر شمارنده

$$R_0 = 1 \quad \blacksquare$$

$$R_1 = 3 \quad \blacksquare$$

$$R_2 = 1 \quad \blacksquare$$

$$n \approx \frac{3}{\phi} 2^{\frac{1}{3} \sum_{i=0}^2 R_i} \approx \frac{3}{0.77351} 2^{\frac{1+3+1}{3}} \approx 12.31$$

تخمین حاصل نزدیک‌تر به مقدار واقعی تعداد منحصر به فرد ۱۰

بدون استفاده از شمارنده‌های زیاد بهتر از مثال قبلی \blacksquare

شمارش احتمالاتی-ویژگی ها

مناسب بودن الگوریتم فلاژوله-مارتین برای مجموعه‌های داده با تعداد منحصر به فرد زیاد

$$\frac{n}{m} > 20 \text{ تقریب خوب در نسبت‌های}$$

جهت بهبود نتایج در مجموعه داده‌های با تعداد منحصر به فردهای کوچک

- امکان استفاده از الگوریتم‌های غیرخطی در الگوریتم
- بیورن شویرمان و مارتین ماوه در سال ۱۳۸۶ پیشنهاد بهبودی احتمالی
- متناسب‌سازی با افزودن عبارتی برای تعداد منحصر به فردهای کوچک
- همگرا شدن به صفر برای تعداد منحصر به فردهای بزرگ

$$n \approx \frac{m}{\phi} (2^{\bar{R}} - 2^{-\chi \times \bar{R}})$$

$$\chi \approx 1.75$$

شمارش احتمالاتی-ویژگی ها

رابطه معکوس خطای استاندارد δ الگوریتم فلاژوله-مارتین با تعداد شمارنده‌های مورد استفاده

$$\delta \approx \frac{0.78}{\sqrt{m}}$$


شمارش احتمالاتی-ویژگی ها

مقادیر مرجع خطای استاندارد برای تعداد شمارنده‌های پرکاربرد
طول ℓ هر شمارنده

$$\ell > \log_2 \left(\frac{n}{m} \right) + 4$$

▪ کفایت $\ell = 32$

شمارش تعداد منحصربه‌فردهای بسیار فراتر از 10^9 با استفاده از ۶۴ شمارنده


$$\begin{aligned}\binom{n}{m} / 2^\ell &\leq 2^{-4} \\ \frac{n}{m} &\leq 2^{\ell-4} \\ \ell - 4 &\geq \log \frac{n}{m} \\ \ell &\geq \log_2 \binom{n}{m} + 4\end{aligned}$$

شمارش احتمالاتی-ویژگی ها

ادغام شمارنده‌های متناظر مجموعه داده‌های مختلف با اعمال عملیات OR بیتی

- نتیجه: شمارنده واحد برای کل مجموعه داده‌ها

عدم پشتیبانی از حذف مانند فیلتر بلوم،

- افزودن امکان حذف با پیروی از رویکرد فیلتر بلوم شمارنده

- پشتیبانی از حذف‌های احتمالا درست

- گسترش آرایه‌های بیتی داخلی با شمارنده‌ها

- منجر به افزایش فضای ذخیره‌سازی

روش لگ لگ

الگوریتم‌های احتمالاتی تخمین تعداد منحصر به فرد پر کاربرد

- با استفاده عملی
- خانواده لگ لگ از الگوریتم‌ها
- شامل الگوریتم لگ لگ پیشنهادی ماریان دوراند و فیلیپ فلاژوله در سال ۱۳۸۲
- دو روش جانشین آن ابر لگ لگ و ابر لگ لگ ++
- استفاده از رویکردی مشابه الگوریتم شمارش احتمالاتی

تخمین تعداد منحصر به فرد n با مشاهده حداکثر تعداد صفرهای ابتدایی در نمایش دودویی مقادیر

نیاز به حافظه کمکی

پیمایش یکباره داده‌ها جهت تولید تخمینی از تعداد منحصر به فرد

روش لگ لگ

مطابق منوال روش‌های قبلی تعیین تعداد منحصر بفرد اجزا

پیش‌پردازش هر مقدار در مجموعه داده با اعمال تابع درهم‌سازی

▪ تبدیل مقادیر به اعداد صحیح با توزیع کافی یکنواخت در بازه $\{0, 1, \dots, 2^{\ell-1}\}$ یا به طور معادل، بر روی مجموعه رشته‌های دودویی به طول ℓ

$$h(x) = i = \sum_{j=0}^{\ell-1} i_j 2^j = (i_{\ell-1} \dots i_1 i_0)_2, i_j \in \{0, 1\}$$

روش لگ لگ

مراحل الگوریتم‌ها مشابه PCSA

ابتدا، تقسیم مجموعه داده اولیه یا جریان ورودی به تعدادی زیرمجموعه

▪ نمایه شدن هر کدام با یکی از m شمارنده

تابع درهم‌ساز واحد

▪ همانند میانگین‌گیری تصادفی

▪ انتخاب شمارنده برای مقدار خاص x

▪ استفاده از بخشی از مقدار درهم $h(x)$

▪ بخش دیگر برای به‌روزرسانی شمارنده مربوطه

روش لگ لگ

همه الگوریتم‌های مورد بحث بر اساس مشاهده الگوهای 10^j

- در ابتدای مقادیر برای شمارنده خاص
- هر الگو با شاخص آن، به نام رتبه،

نمایه شدن رتبه معادل موقعیت ۱ با کمترین اهمیت در نمایش دودویی مقدار درهم

مشاهده تعداد منحصربه‌فرد در هر شمارنده بر اساس رتبه‌های دیده شده

تخمین نهایی تعداد منحصربه‌فرد از چنین مشاهداتی با استفاده از تابع ارزیابی

روش لگ لگ

ذخیره‌سازی

- شمارنده‌ها در الگوریتم شمارش احتمالاتی نسبتاً پرهزینه برای نگهداری
- در مقابل: الگوریتم لگ لگ راه‌حلی کارآتر از نظر ذخیره‌سازی
- همراه با تابع ارزیابی بهتر و رویکرد تصحیح سوگیری

روش لگ لگ

الگوریتم لگ لگ

بر اساس محاسبه رتبه هر مقدار ورودی

با تابع درهم ساز واحد h

فرض n تعداد کل مقادیر نمایه شده در شمارنده

▪ ادعا $\frac{n}{2^j}$ مقدار اختیار کردن رتبه j

▪ پس رتبه بیشینه مشاهده شده نشانه و تقریب مناسبی از مقدار $\log_2 n$

$$R = \max_{x \in D} \left(\text{رتبه}(h(x)) \right) \approx \log_2 n$$

روش لگ لگ

تخمین مذکور دارای خطایی با مرتبه $1,87$ برابر مقدار واقعی
▪ عملاً بی معنی

برای کاهش خطا،

- استفاده از سطل بندی مبتنی بر میانگین گیری تصادفی
- تقسیم مجموعه داده به $m = 2^p$ زیرمجموعه $\{S_0, S_1, \dots, S_{m-1}\}$
- پارامتر دقت p تعداد بیت های استفاده شده در پیمایش

بنابراین، برای هر مقدار x از مجموعه داده، p بیت اول مقدار درهم ℓ -بیتی $h(x)$ یافتن شاخص j زیرمجموعه مناسب

$$j = (i_{p-1} \dots i_1 i_0)_2$$

نمایه شدن بقیه $(\ell - p)$ بیت در شمارنده j -ام

محاسبه رتبه و مشاهده R_j

روش لگ لگ

با توزیع منصفانه،

▪ هر زیرمجموعه دارای n/m عضو

▪ پس، R_j ها در شمارنده‌ها $\{COUNTER[j]\}_{j=0}^{m-1}$ نشانه‌ای از مقدار $\log_2 \frac{n}{m}$

▪ استفاده از میانگین حسابی آنها و تصحیح سوگیری، کاهش وردائی مشاهده

$$n = \alpha_m \cdot m \cdot 2^{\frac{1}{m} \sum_{j=0}^{m-1} R_j}$$

$$\alpha_m = \left(\Gamma\left(-\frac{1}{m}\right) \times \frac{1 - \frac{1}{2^{\frac{1}{m}}}}{\log 2} \right)^m \approx 0.39701 - \frac{2\pi^2 + (\ln 2)^2}{48m}$$

$$\Gamma(\alpha) = \int_0^{\infty} e^{-x} x^{\alpha-1} dx$$

$$\Gamma(n) = (n-1)!$$

برای اکثر موارد عملی $m \geq 64$ کفایت استفاده از $\alpha_m \approx 0.39701$

روش لگ لگ

الگوریتم ۸- تخمین تعداد منحصر به فرد با LogLog

ورودی: مجموعه داده D

ورودی: آرایه‌ای از m شمارنده LogLog با تابع درهم‌ساز h

خروجی: تخمین تعداد منحصر به فرد

$$COUNTER[j] \leftarrow 0, j = 0 \dots m-1$$

برای $x \in D$ انجام بده

$$i \leftarrow h(x) = (i_{\ell-1} \dots i_1 i_0)_2, i_k \in \{0, 1\}$$

$$j \leftarrow (i_{p-1} \dots i_1 i_0)_2$$

$$r \leftarrow \text{rank} \left((i_{\ell-1} \dots i_{p+1} i_p)_2 \right)$$

$$COUNTER[j] \leftarrow \max(COUNTER[j], r)$$

$$R \leftarrow \frac{1}{m} \sum_{j=0}^{m-1} COUNTER[j]$$

برگرداندن $\alpha_m \cdot m \cdot 2^R$

روش لگ-لگ-مثال

فرض کنید $m=4$ پس در نتیجه $p = \log_4=2$ بیت جهت انتخاب شمارنده

خروجی درهم طولی برابر هشت بیت

پس بیت‌های صفر و یک برای شمارنده و بقیه شش بیت برای محاسبه رتبه

<u>(b₇b₆b₅b₄b₃b₂)</u>	<u>(b₁b₀)</u>	$h(x)$ (^ bits)	شهر
۰۱۱۰۱۰	۱۰ = ۲	۰۱۱۰۱۰۱۰	Tehran
۱۰۱۱۰۰	۱۱ = ۳	۱۰۱۱۰۰۱۱	Paris
۰۰۰۱۱۱	۰۰ = ۰	۰۰۰۱۱۱۰۰	London
۱۱۰۰۱۰	۱۰ = ۲	۱۱۰۰۱۰۱۰	Berlin
۰۰۱۰۱۱	۰۱ = ۱	۰۰۱۰۱۱۰۱	Rome
۰۱۰۰۱۱	۱۰ = ۲	۰۱۰۰۱۱۱۰	Madrid
۱۰۰۰۰۱	۱۱ = ۳	۱۰۰۰۰۱۱۱	Kyiv
۱۱۱۰۰۰	۰۰ = ۰	۱۱۱۰۰۰۰۰	Vienna

روش لگ-لگ-مثال

حساب رتبه بیت‌های بالایی

رتبه	<u>(b₇b₆b₅b₄b₃b₂)</u>	شهر
۱	۰۱۱۰۱۰	Tehran
۲	۱۰۱۱۰۰	Paris
۰	۰۰۰۱۱۱	London
۱	۱۱۰۰۱۰	Berlin
۰	۰۰۱۰۱۱	Rome
۰	۰۱۰۰۱۱	Madrid
۰	۱۰۰۰۰۱	Kyiv
۳	۱۱۱۰۰۰	Vienna

روش لگ-لگ-مثال

بروز شدن شمارنده‌ها با عملیات بیش

شمارنده j	شهرهای تخصیص یافته	رتبه	بیش شمارنده
۰	لندن با $r=0$ و وین با $r=3$	$\{0,3\}$	۳
۱	رم با $r=0$	$\{0\}$	۰
۲	تهران با $r=1$ و برلین با $r=1$ و مادرید با $r=0$	$\{0,1,1\}$	۱
۳	پاریس با $r=2$ و کیف با $r=0$	$\{0,2\}$	۲

روش لگ-لگ-مثال

پس شمارنده دارای مقدار $[3,0,1,2]$ است

محاسبه میانگین را به صورتی زیر حساب می‌کنیم.

$$R = \frac{1}{m} \sum [j]_{\text{شمارنده}} = \frac{1}{4} (3 + 0 + 1 + 2) = 1.5$$

نتیجه نهایی تخمین برای $m=4$ و $\alpha_4 \approx 0.39701$ برابر است با

$$\hat{n} = \alpha_m \times m \times 2^R = 0.39701 \times 4 \times 2^{1.5}$$

روش لگ-لگ-ویژگی‌ها

رابطه معکوس خطای استاندارد δ الگوریتم لگ-لگ با تعداد شمارنده‌های استفاده شده m

$$\delta \approx \frac{1.3}{\sqrt{m}}$$

$m = 256$ خطای استاندارد حدود هشت درصد

$m = 1024$ حدود چهار درصد

روش لگ-لگ-ویژگی‌ها

در صورت شمارش تا n

فضای ذخیره‌سازی موردنیاز الگوریتم لگ-لگ حدود $O(\log_2 \log_2 n)$

روش لگ-لگ-ویژگی‌ها

در صورت شمارش تا n

فضای ذخیره‌سازی موردنیاز الگوریتم لگ-لگ حدود $O(\log_2 \log_2 n)$

مقدار دقیق کل فضای مورد نیاز الگوریتم برای شمارش تا n

$$m \cdot \log_2 \log_2 \frac{n}{m} (1 + O(1))$$

روش لگ-لگ-ویژگی‌ها

در صورت شمارش تا n

فضای ذخیره‌سازی موردنیاز الگوریتم لگ-لگ حدود $O(\log_2 \log_2 n)$

مقدار دقیق کل فضای مورد نیاز الگوریتم برای شمارش تا n

$$m \cdot \log_2 \log_2 \frac{n}{m} (1 + O(1))$$

در مقایسه با الگوریتم شمارش احتمالاتی

- نیاز هر شمارنده به ۱۶ یا ۳۲ بیت

- الگوریتم لگ-لگ نیاز به شمارنده‌های بسیار کوچکتر $\{COUNTER[j]\}_{j=0}^{m-1}$

- معمولاً هر کدام پنج بیت

- با وجود مصرف حافظه کمتر الگوریتم لگ-لگ نسبت به الگوریتم شمارش احتمالاتی

- دارای میزان ناچیزی از کاهش دقت نسبت به شمارش احتمالاتی

روش لگ لگ

- نیاز به شمارش تعداد منحصربه‌فرد تا 2^{30}
- یعنی حدود ۱ میلیارد، با دقت حدود چهار درصد
- برای چنین خطای استاندارد، $m = 1024$ سطل مورد نیاز
- هر کدام $\frac{n}{m} = 2^{20}$ مقدار
- تعداد بیت مورد نیاز $\log_2 \log_2 2^{20} \approx 4.32$
- یا حدود ۵ بیت به ازای هر سطل (مقداری کمتر از ۳۲)
- تخمین تعداد منحصربه‌فرد تا حدود 10^9 با خطای استاندارد چهار درصد
 - نیاز الگوریتم به $10^9 \times 5$ سطل ۵ بیتی
 - در مجموع ۶۴۰ بایت

روش ابرلگ لگ

فیلیپ فلاژوله، اریک فیوزی، اولیویه گاندوئه، و فردریک مونیه

▪ سال ۱۳۸۶

▪ پیشنهاد الگوریتم ابرلگ لگ بر اساس الگوریتم لگ لگ

استفاده از تابع درهم ساز ۳۲ بیتی، تابع ارزیابی متفاوت و تصحیحات سوگیری مختلف

مشابه الگوریتم لگ لگ،

▪ استفاده ابرلگ لگ از تصادفی سازی برای تقریب تعداد منحصر به فرد مجموعه داده

▪ طراحی برای مدیریت تعداد منحصر به فرد تا 10^9 با یک تابع درهم ساز ۳۲ بیتی

▪ تقسیم مجموعه داده به $m = 2^p$ زیرمجموعه، با دقت $p \in 4 \dots 16$

روش ابرلگ لگ

تمایز تابع ارزیابی الگوریتم ابرلگ لگ از لگ لگ استاندارد

استفاده الگوریتم لگ لگ اصلی از میانگین حسابی

استفاده ابرلگ لگ از نسخه نرمال شده میانگین هارمونیک

$$\hat{n} \approx \alpha_m \cdot m^2 \cdot \left(\sum_{j=0}^{m-1} 2^{-\text{Coutner}[j]} \right)^{-1}$$

که در آن

$$\alpha_m = \left(m \int_0^{\infty} \left(\log_2 \left(\frac{2+x}{1+x} \right) \right)^m dx \right)^{-1} = \frac{1}{2 \ln 2 \left(1 + \frac{1}{m} (3 \ln 2 - 1) + O(m^{-2}) \right)}$$

شهود استفاده از میانگین هارمونیک

▪ کاهش وردائی به دلیل ویژگی آن در هموار کردن توزیع های احتمال چوله

جدول - α_m برای پرکاربردترین مقادیر m

m	α_m
۲۴	۰.۶۷۳
۲۵	۰.۶۹۷
۲۶	۰.۷۰۹
≥ ۲۷	$\frac{۰.۷۲۱۲ \times m}{m+۱.۰۷۹}$

روش ابرلگ لگ

نیاز به تصحیح تخمین به دلیل خطاهای غیرخطی برای محدوده‌های کوچک و بزرگ

فلاژوله و همکاران در نتیجه آزمایش تجربی

- برای تعداد منحصربه‌فردهای کوچک $n < \frac{5}{2}m$ برای دستیابی به تخمین‌های بهتر،
- تصحیح الگوریتم ابرلگ لگ با شمارش خطی با استفاده از تعدادی از شمارنده‌های COUNTER[j] غیرصفر
- در صورت صفر بودن شمارنده، امکان نتیجه‌گیری خالی بودن آن زیرمجموعه خاص با قطعیت

روش ابرلگ لگ

تصحیحات

$$n = \begin{cases} \text{LinearCounter}, & \hat{n} \leq \frac{5}{2}m \wedge \exists j: \text{Counter}[j] \neq 0 \\ -2^{23} \log\left(1 - \frac{\hat{n}}{2^{32}}\right), & \hat{n} > \frac{1}{30} 2^{32} \\ \hat{n}, & \text{دغا} \end{cases}$$

روش ابرلگ لگ

برای $n = 0$

▪ عدم کفایت تصحیح کافی و برگرداندن همیشه تقریباً $0.7m$

روش ابرلگ لگ

استفاده الگوریتم ابرلگ لگ از تابع درهم ساز ۳۲ بیتی

با نزدیک شدن تعداد منحصربه فرد به $4 \times 10^9 \approx 2^{32}$

رسیدن تابع درهم ساز تقریباً به حد خود و افزایش احتمال تصادم

برای چنین محدوده های بزرگ، تخمین تعداد مقادیر درهم مختلف

- استفاده از آن برای تقریب تعداد منحصربه فرد

- عدم امکان نمایش مقادیر بالاتر و تاثیر بر دقت

روش ابرلگ لگ

تابع درهم‌ساز با نگاشت دامنه به مقادیر l بیت

امکان کدگذاری حداکثر 2^l مقدار مختلف

در صورت نزدیک شدن تعداد منحصر به فرد تخمینی n به چنین حدی
▪ محتمل‌تر شدن تصادم‌های درهم‌ساز

روش ابرلگ لگ

هیچ نشانی دال بر عملکرد بهتر توابع درهم ساز پر کاربرد

- (مانند مورمور ۳، MD5، SHA-1، SHA-256)
- نسبت به سایر روش ها در الگوریتم های ابرلگ لگ یا اصلاحات آن وجود ندارد.

روش ابرلگ لگ

الگوریتم ۹- تخمین تعداد منحصر به فرد با ابرلگ لگ

ورودی: مجموعه داده D

ورودی: آرایه ای از m شمارنده LogLog با تابع درهم ساز h

خروجی: تخمین تعداد منحصر به فرد

$COUNTER[j] \leftarrow 0, j = 0 \dots m-1$

برای $x \in D$ انجام بده

$i \leftarrow h(x) = (i_{r_1} \dots i_1 i_0)_r$

$j \leftarrow (i_{p-1} \dots i_1 i_0)_r$

$r \leftarrow \text{rank}((i_{r_1} \dots i_{p+1} i_p)_r)$

$COUNTER[j] \leftarrow \max(COUNTER[j], r)$

$R \leftarrow \sum_{j=0}^{m-1} r^{COUNTER[j]}$

$\hat{n} = \alpha_m \cdot m^r \cdot \frac{1}{R}$

$n \leftarrow \hat{n}$

اگر $\hat{n} \leq \frac{5}{r} m$ آنگاه

$Z \leftarrow \text{تعداد}_{j=0 \dots m-1} (COUNTER[j] = 0)$

اگر $Z \neq 0$ آنگاه

$n \leftarrow m \cdot \log\left(\frac{m}{Z}\right)$

در غیر این صورت اگر $\hat{n} > \frac{1}{r} 2^{2r}$ آنگاه

$n \leftarrow -2^{2r} \cdot \log\left(1 - \frac{n}{2^{2r}}\right)$

برگرداندن n

روش ابرلگ-لگ

مشابه الگوریتم لگ-لگ، بین خطای استاندارد δ و تعداد شمارنده‌ها m نیاز به سبک-سنگینی

$$\delta \approx \frac{1.04}{\sqrt{m}}$$

عدم رشد خطی حافظه با تعداد مقادیر

- برخلاف الگوریتم‌هایی مانند الگوریتم شمارش خطی
- با تخصیص $(\ell-p)$ بیت برای مقادیر درهم و داشتن $m = 2^p$ شمارنده در مجموع،
- حافظه مورد نیاز برابر با

$$[\log_2(\ell + 1 - p)] \cdot 2^p \text{ bits}$$

استفاده الگوریتم از توابع درهم‌ساز ۳۲ بیتی و دقت $p \in 4 \dots 16$ استفاده می‌کند،

- نیازهای حافظه برای داده‌ساختار ابرلگ-لگ برابر با $5 \cdot 2^p$ بیت

روش ابرلگ لگ

الگوریتم ابرلگ لگ

- تخمین تعداد منحصربه‌فردهای بسیار فراتر از 10^9
- با دقت دو درصد
- با استفاده از تنها ۱,۵ کیلوبایت حافظه

پایگاه داده معروف درون حافظه Redis

- حفظ ساختارهای داده ابرلگ لگ با ۱۲ کیلوبایت
- تقریب تعداد منحصربه‌فردها با خطای استاندارد ۰.۸۱

بهبود تخمین تعداد منحصر به فرد در ابرلگ لگ در مقایسه با لگ لگ در مجموعه‌های داده کوچک

- با این وجود تخمین بیش از حد تعداد منحصربه‌فردهای واقعی

پیاده سازی انواع الگوریتم‌های ابرلگ لگ در پایگاه‌های داده شناخته شده

- مانند Amazon Redshift، Redis، Apache CouchDB، Riak

روش ابرلگ لگ ++

اشتفان هویله، مارک نونکسر، و اسکندر هال

▪ سال ۱۳۹۲

▪ معرفی نسخه بهبود یافته ابرلگ لگ را به نام ابرلگ لگ ++،

▪ متمرکز بر روی تعداد منحصر به فردهای بزرگ و تصحیح سوگیری

استفاده از تابع درهم ساز ۶۴ بیتی

▪ اهم اصلاحات در الگوریتم ابرلگ لگ ++

مقادیر خروجی طولانی تر تابع درهم ساز

▪ امکان کدگذاری مقادیر بیشتر

امکان تخمین تعداد منحصر به فردهای بسیار بزرگتر از 10^9

با نزدیک شدن تعداد منحصر به فرد به $1.8 \times 10^{19} \approx 2^{64}$

▪ مشکل سازی تصادم های درهم برای ابرلگ لگ ++

روش ابرلگ لگ ++

استفاده الگوریتم ابرلگ لگ ++ از تابع ارزیابی لگ لگ

- بهبود تصحیح سوگیری

اجرای مجموعه‌ای از آزمایش‌ها برای اندازه‌گیری سوگیری

- نتایج حال از داده‌های تجربی جمع‌آوری شده

- برای $n \leq 5m$ امکان تصحیح سوگیری الگوریتم ابرلگ لگ اصلی

روش ابرلگ لگ ++

علاوه بر مقاله اصلی، هویله و همکاران

ذخیره مقادیر تعیین شده تجربی

- در آرایه‌هایی از تخمین‌های خام تعداد منحصربه‌فرد
RAWESTIMATEDATA

- عرضه سوگیری‌های مرتبط BIASDATA برای بهبود تصحیح سوگیری در الگوریتم
سعی پوشش هر مورد ممکن

- RAWESTIMATEDATA عرضه آرایه‌ای از ۲۰۰ نقطه درون‌یابی

- حاوی میانگین تخمین خام اندازه‌گیری شده در نقاط در بیش از ۵۰۰۰ مجموعه داده مختلف
BIASDATA

- شامل حدود ۲۰۰ سوگیری اندازه‌گیری شده و مطابق با RAWESTIMATEDATA

روش ابرلگ لگ ++

- هر دو آرایه صفر-اندیس
- حاوی مقادیر از پیش محاسبه شده برای همه دقت‌های پشتیبانی شده $p \in 4 \dots 18$
- تطابق اندیس صفر در آرایه‌ها با مقدار دقت ۴
- مثال، برای $m = 2^{10}$ و $p = 10$
- یافتن داده‌های مورد نیاز در `BIASDATA[6]` و `RAWESTIMATEDATA[6]`

روش ابرلگ-لگ

الگوریتم ۱۰ - تصحیح سوگیری با استفاده از مقادیر تجربی

ورودی: تخمین \hat{n} با دقت p

خروجی: تخمین تعداد منحصر به فرد تصحیح شده با سوگیری

$n_{low} \leftarrow \cdot, n_{up} \leftarrow \cdot, j_{low} \leftarrow \cdot, j_{up} \leftarrow \cdot$

برای $j \leftarrow \cdot$ تا طول $RAWESTIMATEDATA[p - 4]$ انجام بده

اگر $RAWESTIMATEDATA[p - 4][j] \geq \hat{n}$ آنگاه

$j_{low} \leftarrow j - 1, j_{up} \leftarrow j$

$n_{low} \leftarrow RAWESTIMATEDATA[p - 4][j_{low}]$

$n_{up} \leftarrow RAWESTIMATEDATA[p - 4][j_{up}]$

$break$ |

$b_{low} \leftarrow BIASDATA[p - 4][j_{low}]$

$b_{up} \leftarrow BIASDATA[p - 4][j_{up}]$

$y =$ درونیابی $((n_{low}, n_{low} - b_{low}), (n_{up}, n_{up} - b_{up}))$

برگرداندن $y(\hat{n})$

روش ابرلگ لگ ++ -- مثال

تصحیح سوگیری با استفاده از مقادیر تجربی

فرض: حاصل تخمین تعداد منحصر به فرد $\hat{n} = 2018.34$ از روش اصلی

در پی تصحیح آن برای دقت $p = 10$
($m = 2^{10}$) .

روش ابرلگ لگ ++

بررسی آرایه $RAWESTIMATEDATA[6]$

- تعیین اندیس‌های دو سر بازه حاوی مقدار \hat{n}
- بین مقادیر با اندیس‌های ۷۳ و ۷۴ آرایه

▪ $RAWESTIMATEDATA[6][73] = 2003.1804$ و $RAWESTIMATEDATA[6][74] = 2026.071$

$$2003.1804 \leq \hat{n} \leq 2026.071$$

بازیابی سوگیری‌ها از $BIASDATA[6]$ در موقعیت‌های ۷۳ و ۷۴

▪ $BIASDATA[6][73] = 134.1804$ و $BIASDATA[6][74] = 131.071$

تخمین صحیح در بازه

▪ $[2003.1804 - 134.1804, 2026.071 - 131.071] = [1869.0, 1895.0]$

روش ابرلگ لگ ++

درون یابی برای محاسبه تقریب تصحیح شده

- استفاده از جستجوی k -نزدیک ترین همسایه یا صرفاً با درون یابی خطی $y(x) = ax + b$
- $y(2003.1804) = 1869.0$ و $y(2026.071) = 1895.0$

حاصل درون یابی

$$y = 1.135837 x - 406.28725$$

مقدار درون یابی شده برای تخمین تعداد منحصر به فرد

$$n = y(\hat{n}) = y(2018.34) \approx 1886.218$$

روش ابرلگ لگ ++

با توجه به آزمایش‌های نویسندگان ++ *HyperLogLog*.

▪ بهتر بودن تخمین n_{lin} بر اساس الگوریتم شمارش خطی در مقایسه با مقدار تصحیحی با سوگیری n برای تعداد منحصر به فردهای کوچک

در صورت وجود حداقل یک شمارنده خالی

▪ محاسبه تخمین خطی و فهرستی از آستانه‌های تجربی (در جدول ادامه)

استفاده از مقدار تصحیح شده با سوگیری n

▪ تخمین خطی n_{lin} بالاتر از آستانه χ_m برای m فعلی

روش ابرلگ لگ ++ -- مثال

تصحیح سوگیری با آستانه

$$n \approx 1886.218 \quad m = 2^{10} \text{ محاسبه مقدار تصحیح شده با سوگیری}$$

جهت تعیین ترجیح این مقدار به تخمین شمارش خطی

▪ نیاز به یافتن تعداد شمارنده‌های خالی Z

▪ فرض $Z = 73$

$$n_{lin} = 2^{10} \cdot \log\left(\frac{2^{10}}{73}\right) \approx 2704 \text{ تخمین خطی}$$

مقایسه n_{lin} با آستانه $\chi_m = 900$ از جدول

▪ بسیار کمتر از مقدار محاسبه شده

▪ ترجیح تخمین تصحیح شده با سوگیری n به تخمین شمارش خطی n_{lin}

روش ابرلگ لگ ++

آستانه‌های تجربی χ_m برای مقادیر دقت پشتیبانی شده

p	m	χ_m	p	m	χ_m	p	m	χ_m
۴	2^4	۱۰	۹	2^9	۴۰۰	۱۴	2^{14}	۱۱۵۰۰
۵	2^5	۲۰	۱۰	2^{10}	۹۰۰	۱۵	2^{15}	۲۰۰۰۰
۶	2^6	۴۰	۱۱	2^{11}	۱۸۰۰	۱۶	2^{16}	۵۰۰۰۰
۷	2^7	۸۰	۱۲	2^{12}	۳۱۰۰	۱۷	2^{17}	۱۲۰۰۰۰
۸	2^8	۲۲۰	۱۳	2^{13}	۶۵۰۰	۱۸	2^{18}	۳۵۰۰۰۰

روش ابرلگ لگ ++

الگوریتم ۱۱- تخمین تعداد منحصر به فرد با ابرلگ لگ ++

ورودی: مجموعه داده D

ورودی: آرایه ای از m شمارنده لگ لگ با تابع درهم ساز h

خروجی: تخمین تعداد منحصر به فرد

$COUNTER[j] \leftarrow 0, j = 0 \dots m-1$

برای $x \in D$ انجام بده

$i \leftarrow h(x) = (i_{p-1} \dots i_1 i_0)_r$

$j \leftarrow (i_{p-1} \dots i_1 i_0)_r$

$r \leftarrow rank((i_{p-1} \dots i_{p+1} i_p)_r)$

$COUNTER[j] \leftarrow max(COUNTER[j], r)$

$R \leftarrow \sum_{k=0}^{m-1} r^{COUNTER[k]}$

$\hat{n} = \alpha_m \cdot m^2 \cdot \frac{1}{R}$

$n \leftarrow \hat{n}$

اگر $\hat{n} \leq 5m$ آنگاه

$n \leftarrow CorrectBias(\hat{n})$

$Z \leftarrow \sum_{j=0}^{m-1} (COUNTER[j] = 0)$

اگر $Z \neq 0$ آنگاه

$n_{lin} \leftarrow m \cdot \log \frac{m}{Z}$

اگر $n_{lin} \leq \chi_m$ آنگاه

$n \leftarrow n_{lin}$

برگرداندن n

الگوریتم ۹- تخمین تعداد منحصر به فرد با ابرلگ لگ ++

ورودی: مجموعه داده D

ورودی: آرایه ای از m شمارنده $LogLog$ با تابع

خروجی: تخمین تعداد منحصر به فرد

$COUNTER[j] \leftarrow 0, j = 0 \dots m-1$

برای $x \in D$ انجام بده

$i \leftarrow h(x) = (i_{p-1} \dots i_1 i_0)_r$

$j \leftarrow (i_{p-1} \dots i_1 i_0)_r$

$r \leftarrow rank((i_{p-1} \dots i_{p+1} i_p)_r)$

$COUNTER[j] \leftarrow max(COUNTER[j], r)$

$R \leftarrow \sum_{j=0}^{m-1} r^{COUNTER[j]}$

$\hat{n} = \alpha_m \cdot m^2 \cdot \frac{1}{R}$

$n \leftarrow \hat{n}$

اگر $\hat{n} \leq \frac{5}{r} m$ آنگاه

$Z \leftarrow \sum_{j=0}^{m-1} (COUNTER[j] = 0)$

اگر $Z \neq 0$ آنگاه

$n \leftarrow m \cdot \log \left(\frac{m}{Z} \right)$

در غیر این صورت اگر $\hat{n} > \frac{1}{r} 2^{32}$ آنگاه

$n \leftarrow -2^{32} \cdot \log \left(1 - \frac{n}{2^{32}} \right)$

برگرداندن n

روش ابرلگ لگ ++

دقت بهتر یا متناسب ابرلگ لگ ++ با ابرلگ لگ

تعداد منحصربه فردهای بین ۱۲۰۰۰ و ۶۱۰۰۰،

▪ امکان خطای کمتر با تصحیح سوگیری

▪ اجتناب از افزایش ناگهانی خطا هنگام جابجایی بین الگوریتم‌های فرعی

روش ابرلگ لگ ++

عدم نیاز ابرلگ لگ به ذخیره مقادیر درهم،

- فقط یک به اضافه حداکثر اندازه تعداد صفرهای ابتدایی
- عدم افزایش قابل توجه حافظه مورد نیاز در مقایسه با ابرلگ لگ
- با توجه به جدول، نیاز به $6 \cdot 2^p$ بیت

تخمین تعداد منحصر به فرد حدود $10^9 * 7.9$ با میزان خطای ۱.۶۲۵ درصد
▪ با استفاده از ۲,۵۶ کیلوبایت حافظه

روش ابرلگ لگ ++

استفاده الگوریتم از رویکرد میانگین گیری تصادفی

تقسیم مجموعه داده به $m = 2^p$ زیرمجموعه $\{S_j\}_{j=0}^{m-1}$

هر کدام دارای شمارنده‌های مرتبط $\{COUNTER[j]\}_{j=0}^{m-1}$

هر شمارنده مدیریت اطلاعات مربوط به $\frac{n}{m}$ مقدار

عدم استفاده از بیشتر شمارنده‌ها وقتی $n \ll m$

- عدم نیاز به ذخیره‌سازی
- امکان استفاده از نمایش تنک
- در صورت کوچک‌تر بودن تعداد منحصر به فرد n از m
- نیاز به حافظه بسیار کمتر ابرلگ لگ ++ نسبت به روش‌های پیشین

الگوریتم ابرلگ‌لگ ++ نسخه تنک

صرفاً ذخیره جفت‌های $(j, COUNTER[j])$

نمایش آنها با الحاق الگوهای بیتی خود به عنوان یک عدد صحیح منفرد

ذخیره همه جفت‌ها در فهرست مرتب شده منفرد از اعداد صحیح

همیشه در پی محاسبه حداکثر رتبه، عدم نیاز به ذخیره جفت‌های مختلف با اندیس یکسان
▪ در عوض ذخیره جفتی با حداکثر اندیس

در عمل، برای ارائه تجربه بهتر، امکان ایجاد فهرست مرتب نشده دیگر برای درج‌های سریع

مرتب‌سازی دوره‌ای و ادغام در فهرست اصلی

ذیل!

سوراو چاکرابورتی، وینودچاندران، کولدیپ میل
۱۴۰۱

معرفی الگوریتمی برای مسئله زیر

▪ ورود یک به یک دنباله‌ای از مقادیر (a_1, a_2, \dots, a_m)

▪ به دنبال یافتن تعداد متمایز آنها

▪ $A = \{a_1, a_2, \dots, a_m\}$ مجموعه عناصر موجود در این دنباله

▪ با چشمپوشی از تکرارها

▪ به دنبال اندازه مجموعه یعنی $|A|$

▪ $|A|$ احتمالاً بسیار بزرگتر از تعداد مقادیر قابل‌نگهداری در حافظه در هر لحظه

راهبرد مناسب محاسبه تخمین ناریب از $|A|$

ذیل

الگوریتمی جالب و در عین حال بسیار ساده

به پیشنهاد کنوت در سال ۱۴۰۲ معروف به الگوریتم CVM

الگوریتم دارای بافری با اندازه محدود

▪ جهت حفظ زیرمجموعه‌ای تصادفی از آن مقادیر

▪ $\Pr[(1 - \epsilon) \cdot |A| \leq c \leq (1 + \epsilon) \cdot |A|] \geq 1 - \delta$

مقدار ورودی	$h(x)$ (عدد صحیح 0-255)
A	120
B	200
C	45
D	200
E	12
F	240
G	78

A, B, C, D, A, E, B, C, F, G
 $n=7$ و $\{A, B, C, D, E, F, G\}$

روش‌های شمارش خطی، PCSA، بهبود شویرمان، لگ‌لگ و ابرلگ‌لگ
شمارش خطی با $m=16$
PCSA و بقیه موارد با چهار شمارنده
میزان خطای هر یک را حساب کنید

مقدار ورودی	$h(x)$ (عدد صحیح 0-255)
A	120
B	200
C	45
D	200 (تصادم با B)
E	12
F	240
G	78

A, B, C, D, A, E, B, C, F, G
 $n=7$ و $\{A, B, C, D, E, F, G\}$

روش‌های شمارش خطی، PCSA، بهبود شویرمان، لگ‌لگ و ابرلگ‌لگ

شمارش خطی با $m=16$

PCSA و بقیه موارد با چهار شمارنده

میزان خطای هر یک را حساب کنید

مقدار ورودی	$h(x)$ (عدد صحیح 0-255)	$h(x)$ (دو مقداری هشت بیتی)
A	120	01111000
B	200	11001000
C	45	00101101
D	200 (تصادم با B)	11001000
E	12	00001100
F	240	11110000
G	78	01001110

A, B, C, D, A, E, B, C, F, G
 $n=7$ و $\{A, B, C, D, E, F, G\}$

روش‌های شمارش خطی، PCSA، بهبود شویرمان، لگ‌لگ و ابرلگ‌لگ

شمارش خطی با $m=16$

PCSA و بقیه موارد با چهار شمارنده

میزان خطای هر یک را حساب کنید

```
class CVM:
    def __init__(self, epsilon, delta):
        self.s = ceil((24 / epsilon**2) * log(8 / delta))
        self.buffer = set()
        self.threshold = 1.0

    def add(self, x):
        self.buffer.add(x)

        while len(self.buffer) > self.s:
            # Purge each element with probability 0.5
            survivors = {e for e in self.buffer
                        if random.random() < 0.5}
            self.buffer = survivors
            self.threshold *= 2

    def estimate(self):
        return len(self.buffer) * self.threshold
```