

تعداد متمایز در مجموعه داده با مقادیر تکراری

محسن هوشمند

فهرست

۲	تعداد منحصر به فرد
۳	شمارش ارقام متمایز در پایگاه‌های داده
۴	شمارش خطی
۸	شمارش احتمالاتی
۱۳	ویژگی‌ها
۱۴	روش‌های لگ‌لگ و ابرلگ‌لگ
۱۴	الگوریتم لگ‌لگ
۱۶	ویژگی‌ها
۱۷	الگوریتم ابرلگ‌لگ
۱۹	ویژگی‌ها
۲۰	الگوریتم ابرلگ‌لگ ++
۲۲	ویژگی‌ها

تعداد منحصر به فرد

مسئله تخمین تعداد منحصر به فرد در پی یافتن تعداد مقادیر متمایز در مجموعه داده‌ای دارای اعضای تکراری است که در حوزه‌هایی گوناگون تولید نرم‌افزار اعم از پایگاه داده‌ها و ترافیک شبکه کاربرد دارد. در روش سنتی و دم دستی تعیین تعداد دقیق یک مجموعه، فهرستی از همه مقادیر ایجاد، و سپس با مرتب‌سازی و جستجو از شمردن عضو تکراری اجتناب و در نتیجه تعداد مقادیر متمایز احصاء می‌شود. هنگامی که تعداد کل اعضای مجموعه از جمله تکراری‌ها برابر N باشد، با روش مذکور پیچیدگی زمانی شمارش دقیق مقادیر منحصر به فرد برابر با $O(N \log N)$ است و نیاز به حافظه کمکی الزامی است. اما با گسترش خدمات اینترنتی، با میلیاردها کلیک، جستجو، و سفارش که روزانه با تعداد بسیار کمتری کاربر انجام می‌پذیرد، بازطراحی شمارش تعداد متمایز را اهمیتی دوباره بخشیده است. خاصه، تعیین تعداد منحصر به فرد با پیمایش یکباره داده و در فضای اشغالی کمتر از حوزه‌های مهم تحقیقاتی به شمار می‌رود.

در حال حاضر، تخمین کاردینالیته برای تعیین تعداد بازدیدکننده یکتا از محصولی خاص، تعداد کاربران متمایز استفاده کننده از قابلیت مشخص در وب‌اپلیکیشن، و همچنین برای تشخیص تغییرات ناگهانی در تعداد آدرس‌های IP مبدأ-مقصد یکتا که از یک روتر عبور می‌کنند (که ممکن است نشان‌دهنده حمله انکار خدمت باشد) به کار می‌رود. به علاوه، به دلیل نحوه تکثیر مکرر اطلاعات در وب، اندازه‌گیری کاردینالیته در تشخیص تعداد محتوای متمایز برای مثال، تعداد مقاله‌های خبری یکتا یا تعداد نسخه‌های متفاوت از محتوای وبسایت یاری‌رسان است. نظر به مجموعه داده‌های عظیم امروزی، علاقه فزاینده‌ای به طراحی الگوریتم‌هایی وجود دارد که بتوانند تعداد منحصر به فرد مجموعه را با دقتی مناسب و در فضایی بسیار کوچک‌تر از اندازه خود مجموعه تقریب بزنند.

مثال ۱- بازدیدکنندگان منحصر به فرد

یکی از شاخص‌های کلیدی ارزشمند KPIs برای هر تارمانه، تعداد بازدیدکنندگان منحصر به فرد آن در بازه زمانی مشخص است. جهت سادگی، فرض می‌کنیم که هر بازدیدکننده از یک نشانی IP در تمامی بازدیدهایش استفاده می‌کند، در نتیجه، با شمارش تعداد IP-های منحصر به فرد که با پروتکل اینترنت نسخه IP_v6 با رشته‌های 128 بیتی نشان داده می‌شوند تعداد بازدیدکنندگان احصاء می‌شود. آیا این کار آسانی است؟ آیا می‌توان با تکیه صرف بر روش‌های کلاسیک برای شمارش دقیق تعداد را بدست آورد؟ پاسخ به این سوال به میزان استفاده از تارمانه بستگی دارد.

مطابق [SimilarWeb](#) در گزارش آمار ترافیک مارس ۲۰۱۷ (اسفند ۱۳۹۵) سه تارمانه خرده فروشی پر مشتری در ممالک اتا زونی اعم از آمازون و آئی‌بی و والمارت، میانگین تعداد بازدید از این تارمانه‌ها حدود ۱.۴۴ میلیارد و میانگین تعداد صفحات مشاهده شده در هر بازدید ۸.۲۴ بود. سخن کوتاه، آمار اسفند ۱۳۹۵ مراجعه به حدود ۱۲ میلیارد آدرس IP 128 بیتی است و نگهداری همه آنها اندازه‌ای در حدود ۱۹۲ گیگابایت لازم دارد. اگر فرض کنیم که از هر ۱۰ بازدیدکننده، یک نفر منحصر به فرد بوده است، می‌توانیم انتظار داشته باشیم که تعداد منحصر به فرد چنین مجموعه‌ای حدود ۱۴۴ میلیون باشد و حافظه مورد نیاز برای ذخیره فهرست مقادیر منحصر به فرد ۲۳ گیگابایت است.

مثال زیر چالش تخمین تعداد منحصر به فرد را برای تحقیقات علمی نشان می‌دهد.

مثال ۲- تحلیل DNA

مطالعه همبستگی‌ها در توالی‌های DNA، از پژوهش‌های مهم در تحقیقات ژنوم انسانی است. مولکول‌های DNA شامل دو رشته جفت شده هستند که هر کدام از چهار واحد شیمیایی پایه DNA، با علامت A (آدنین)، G (گوانین)، C (سیتوزین) و T (تیمین) تشکیل شده‌اند. ژنوم انسان حدود سه میلیارد جفت پایه دارد. توالی‌یابی به معنای تعیین ترتیب دقیق جفت‌های پایه در قطعاتی از

DNA است. از دید ریاضی، توالی DNA را می‌توان به عنوان رشته‌ای با اندازه‌های طولانی از نمادهای A, G, C, T در نظر گرفت، به طوری که می‌توان آنها را به عنوان نمونه‌ای از مجموعه داده‌ای بالقوه بی‌نهایت در نظر گرفت. مسئله‌اندازه‌گیری همبستگی در پی تعیین تعداد زیررشته‌های متمایز با اندازه ثابت در قطعه‌ای DNA تعریف می‌شود. مفروض مسئله مذکور این است که هر توالی با زیررشته‌های متمایز کم، همبستگی بیشتری نسبت به توالی دیگری با اندازه مشابه اما با زیررشته‌های متمایز بیشتر دارد. چنین آزمایش‌هایی نیاز به اجرای چندباره روی فایل‌های بزرگ متعدد دارند و برای تسریع تحقیقات به حافظه محدود یا حتی ثابت و زمان اجرای کم دسترسی و تاکید دارند، که امر مذکور با الگوریتم‌های شمارش دقیق ممکن نیست.

بنابراین، تخمین دقیق تعداد منحصره‌فرد به زمان پردازش و حافظه زیاد نیاز دارند. برنامه‌های داده بزرگ به سمت رویکردهای عملی‌تر غالباً بر مبنای الگوریتم‌های احتمالاتی و عرضه‌پاسخی تقریبی گرایش دارند.

در خور ذکر است که هنگام پردازش داده‌ها، داشتن تحلیلی از اندازه مجموعه داده ورودی و تعداد مقادیر متمایز احتمالی در تصمیم‌گیری نقشی پررنگ دارد. دنباله بی‌نهایت (بالقوه بی‌نهایت) رشته‌های تک حرفی a, d, s, ... را در نظر می‌گیریم که بر اساس حروف الفبای لاتین باشد. تعداد منحصره‌فرد را می‌توان به راحتی تخمین زد و کران بالای آن به تعداد حروف لاتین محدود می‌شود. بدیهی است که در این مورد، نیازی به استفاده از هیچ رویکرد احتمالاتی نیست و راه‌حل مبتنی بر لغت‌نامه ساده برای محاسبه دقیق تعداد منحصره‌فرد بسیار خوب عمل می‌کند.

در حل مسئله تعداد منحصره‌فرد، بسیاری از روش‌های احتمالاتی مطرح از الگوریتم فیلتر بلوم تاثیر گرفته و بر اساس مقادیر درهم عمل کرده، سپس، الگوهای رایج در توزیع آنها را مشاهده می‌کنند و «حدس‌های» منطقی در مورد تعداد مقادیر منحصره‌فرد بدون نیاز به ذخیره همه آنها می‌زنند.

شمارش اقلام متمایز در پایگاه‌های داده

احتمالاً آشناترین مثال برای سنجش کاردینالیته (اندازه مجموعه) از پایگاه‌های داده و نحوه استفاده SQL از کلیدواژه DISTINCT ناشی می‌شود. وقتی این کلیدواژه بر روی ستونی در جدول اعمال شود، دستور SELECT DISTINCT تمام اقلام متمایز آن ستون را و دستور SELECT COUNT DISTINCT تعداد اقلام متمایز موجود در ستون مورد نظر را برمی‌گرداند. پرس‌وجوهای حاوی COUNT DISTINCT در بسیاری از حوزه‌ها به‌ویژه در حوزه تجارت الکترونیک یا زمانی که آمار استفاده از تارمانه موردنیاز باشد بسیار شایع و رایج است. داده‌های بازدید کاربران اغلب در جدولی به نام «بازدیدهای-روزانه» ثبت می‌شوند که معمولاً بسیار بزرگ است و دارای ویژگی‌هایی نظیر session_id, timestamp, product_id, user_ip_address, visit_duration و غیره است. با اجرای دستور انتخاب زیر:

```
SELECT COUNT (DISTINCT user_ip_address) WHERE product_id = ۹۸۷۳۹۴۷  
FROM DAILY_VISITS
```

تعداد آدرس‌های IP متمایز (یعنی کاربران) دسترسی‌یافته به محصول با شناسه ۹۸۷۳۹۴۷ را در روزی معین دریافت خواهیم کرد. در تارمانه شلوغ، جدول بازدیدهای روزانه می‌تواند تا چند میلیارد ردیف رشد کند، و این پرس‌وجوی خاص ممکن است نیاز به زمان زیادی برای پاسخ داشته و با تأخیر همراه باشد.

این تأخیر عمدتاً ناشی از عملیات مرتب‌سازی است که COUNT DISTINCT کلاسیک در بیشتر پایگاه‌های داده (مانند Azure SQL/SQL Server) اعمال می‌کنند، مگر آنکه ستون از قبل مرتب شده باشد. پس از مرتب‌سازی ستون، تمام رکوردهای تکراری در کنار یکدیگر قرار می‌گیرند، و یک پیمایش ترتیبی کافی است تا اقلام متمایز شناسایی و شمارش شوند. هزینه عملیات مرتب‌سازی در جدولی با n ردیف، $O(n \log n)$ است و نه تنها برای چند میلیارد ردیف بلکه حتی برای چند میلیون ردیف به خوبی مقیاس‌پذیر نیست، حتی پرس‌وجوهای ساده نیز COUNT DISTINCT‌ها و GROUP BY‌های متعددی روی ستون‌های مختلف انجام می‌دهند، و مرتب‌سازی یک ستون به کاهش پیچیدگی مرتب‌سازی ستون دیگر کمکی نمی‌کند. جهت تسریع کار می‌توان از جدول درهم استفاده کرد، اما جدول درهم همچنان به فضایی خطی نسبت به تعداد مقادیر متمایز k نیاز دارد. از آنجا که k می‌تواند تا n افزایش یابد، استفاده از درهم‌سازی نیز مقرون به صرفه نیست.

جهت رفع مسائل مقیاس‌پذیری، نسخه‌های جدیدتر سیستم‌های مدیریت پایگاه‌داده و انبارهای داده به تخمین تعداد منحصر به فرد روی آورده‌اند. محیط ۲۰۱۹ *SQL Server* دارای عملیات `APPROX_COUNT_DISTINCT` (<http://mng.bz/QWjm>) است که از فضای بسیار کمی استفاده و سریع کار می‌کند. Google BigQuery این رویکرد تقریبی و احتمالاتی را در `COUNT DISTINCT` به عنوان پیش‌فرض قرار می‌دهد، و `EXACT_COUNT_DISTINCT` را برای موقعیت‌هایی که مطلقاً به پاسخ دقیق نیاز است، محفوظ نگه می‌دارد (<http://mng.bz/y>). الگوریتمی که در چنین تخمین‌گرهایی اجرا می‌شود ابرلگ نام دارد و آن را فلاژوله و همکاران [۲] ابداع کردند، به طوری که صرفه‌جویی سخت‌زیادی در فضا (در حد کیلوبایت) ارائه می‌دهد، در حالی که مجموعه‌داده‌هایی با تعداد اعضای در حد تریلیون را پردازش می‌کند و میزان خطا را نسبتاً پایین در مرتبه‌ای $O(\frac{1}{\sqrt{m}})$ نگه می‌دارد که m نشان‌دهنده تعداد مدخل‌های حافظه چندبیتی است. انتخاب رایج برای m مقدار 2^{14} است.

نمونه‌های متعددی از صرفه‌جویی در فضا به بهای کاهش دقت معرفی شده‌اند. با این حال، ابرلگ معنای تازه‌ای به کارایی فضا می‌بخشد، به گونه‌ای که تقریباً همیشه در محدوده چند کیلوبایت باقی می‌ماند و در عین حال تخمینی از تعداد منحصر به فرد واقعی را با میانگین خطای کوچک (مثلاً $\pm 2\%$ درصد) به دست می‌دهد. در ادامه، چند روش تخمین تعداد منحصر به فرد را معرفی می‌کنیم.

شمارش خطی

الگوریتم شمارش خطی به عنوان اولین رویکرد احتمالاتی برای مسئله تعداد منحصر به فرد مطرح شد. طرح اصلی را مورتون آسترهان، ماریو شولنیک، و کیو-یانگ وانگ در سال ۱۳۶۶ معرفی کردند و الگوریتم عملی را کیو-یانگ وانگ، براد واندر-زاندن، و هوارد تیلور در سال ۱۳۶۹ پیشنهاد دادند. طرح اصلی پیشنهادی آنها استفاده از درهم مقادیر تابعی درهم‌ساز است که موجب برتری نسبت به روش‌های دقیق شده است. استفاده از درهم امکان حذف تکراری‌ها بدون نیاز به مرتب‌سازی مقادیر را فراهم می‌کند. البته این فرصت بدون هزینه نیست و از دقت پاسخ می‌کاهد و احتمال خطا به دلیل تصادم‌های احتمالی درهم (نمی‌توان تکراری‌ها خاصه «تکراری‌های تصادفی» را تشخیص داد) را موجب می‌شود. سخن کوتاه، استفاده از چنین جدول درهمی به همراه روش پیمایش مناسب منجر به بهبود عملکرد فضا و زمان نسبت به روش کامل یا جستجو با کمک درهم می‌شود.

با این حال، برای مجموعه‌های داده با تعداد منحصر به فرد عظیم، چنین جدول‌های درهمی می‌توانند بسیار بزرگ باشند و به حافظه‌ای نیاز دارند که به صورت خطی با تعداد مقادیر متمایز در مجموعه رشد می‌کند. برای سیستم‌هایی با حافظه محدود، در نقطه‌ای به ذخیره‌سازی دیسک یا توزیع شده نیاز خواهد داشت، که به دلیل دسترسی کند به دیسک یا شبکه، مزایای جدول‌های درهم سخت کاهش می‌یابد.

مشابه ایده فیلتر بلوم، برای دور زدن چنین مشکلی، الگوریتم شمارش خطی مقادیر درهم را ذخیره نمی‌کند، بلکه بیت‌های متناظر آنها را ذخیره می‌کند و آرایهٔ بیتی (تحت نام «شمارندهٔ خطی») را جانشین جدول درهم می‌کند. فرض می‌شود که مقدار m با تعداد مقادیر متمایز مورد انتظار n متناسب است، اما برای اکثر موارد فقط به یک بیت به ازای هر مقدار نیاز دارد که قابل اجرا است. در ابتدا، همه بیت‌ها در «شمارندهٔ خطی» برابر با صفر هستند. با افزودن مقدار جدید x به چنین داده‌ساختاری، مقدار درهم‌ساز $h(x)$ متناظرش محاسبه و مقدار بیت در نمایهٔ متناظر در شمارنده برابر یک می‌شود.

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

$LINEARCOUNTER[j] \leftarrow 1$

به دلیل استفاده از تک تابع درهم‌ساز h ، امکان دارد دو مقدار ورودی متفاوت منجر به مقداردهی یک بیت در آرایه می‌شود و در نتیجه به تعداد تصادم‌های فراوانی منجر شود. بنابراین، تعداد دقیق (یا حتی نزدیک به دقیق) مقادیر متمایز دیگر نمی‌تواند مستقیماً از چنین طرح فشرده‌ای به دست آید.

الگوریتم منجر به توزیع مقادیر در سطل‌ها (بیت‌های نمایه شدهٔ مقادیر درهم) می‌شود و آرایه بیتی نمایشگر موجب می‌شود مدخل‌های متناظر مقداردهی شوند. مشاهده تعداد یک‌ها در آرایه تخمین تعداد منحصربه‌فرد را به دست می‌دهد.

در مرحله اول الگوریتم شمارش خطی، داده‌ساختار را همانطور که در الگوریتم ۱ نشان داده شد، ایجاد کمی‌کند. با داشتن چنین طرح فشرده‌ای، تعداد منحصربه‌فرد را می‌توان با استفاده از کسر مشاهده شده بیت‌های خالی V با معادله تخمین زد:

$$n \approx -m \cdot \ln V \quad (1)$$

تصادم‌ها بر تخمین تعداد منحصربه‌فرد در الگوریتم شمارش خطی تأثیر می‌گذارند. به بیان دیگر، هر تصادم تعداد بیت‌هایی که باید تنظیم شوند را کاهش می‌دهد و کسر مشاهده شده بیت‌های تنظیم نشده را بزرگتر از مقدار واقعی می‌کند. اگر تصادم درهم وجود نداشت، تعداد نهایی بیت‌های تنظیم شده، تعداد منحصربه‌فرد مورد نظر بود. با این حال، تصادم‌ها اجتناب‌ناپذیر هستند و معادله (۱) در واقع تخمین بیش از حدی از تعداد دقیق منحصربه‌فرد ارائه می‌دهد و از آنجایی که تعداد منحصربه‌فرد مقداری عدد صحیح است، ترجیح می‌دهیم نتیجه آن را به نزدیک‌ترین عدد صحیح کوچکتر گرد کنیم. شبیه محاسبات یافتن تعداد منحصر به فرد در الگوریتم بلوم می‌توان مقدار تخمین معادله (۱) را به صورت زیر به دست آورد. 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 تعدد بیت‌های متناظر یک در فیلتر است. جهت تقریب n از N می‌توان از $N/m = 1 - e^{-n/m}$ مقدار زیر را بدست آورد:

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

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

لگاریتم طبیعی دو طرف را محاسبه می کنیم

$$-\frac{n}{m} = \ln\left(1 - \frac{N}{m}\right)$$

با قرار دادن $V = 1 - \frac{N}{M}$ داریم

$$n = -m \ln V$$

که دقیقا برابر مقداری است که الگوریتم وقتی $0 < N < m$ برمی گرداند. سخن کوتاه، الگوریتم کامل شمارش را به صورت زیر تدیون می کنیم.

الگوریتم ۲- تخمین تعداد منحصره فرد با شمارش خطی

ورودی: مجموعه داده D

خروجی: تخمین تعداد منحصره فرد

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

برای هر $x \in D$ انجام بده

$Add(x, LINEARCOUNTER)$

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

برگرداندن $\lfloor -m \times \ln \frac{Z}{m} \rfloor$

مثال ۳- الگوریتم شمارش خطی: مجموعه داده ای شامل بیست نام از پایتختها برلین، برلین، پاریس، برلین، لیسبون، کیف، پاریس، لندن، رم، آتن، مادرید، وین، رم، رم، لیسبون، برلین، پاریس، لندن، کیف، واشنگتن را در نظر می گیریم. برای چنین تعداد منحصره فرد کوچک (تعداد دقیق منحصره فرد ۱۰ است) برای داشتن خطای استاندارد حدود ده درصد نیاز است طول داده ساختار را حداقل به اندازه تعداد مقادیر منحصره فرد مورد انتظار انتخاب کنیم، بنابراین $m = 2^4$ را انتخاب می کنیم. به عنوان تابع درهم ساز h با مقادیر در $\{0, 1, \dots, 2^4 - 1\}$ از تابعی مبتنی بر مورمور ۳ نوع ۳۲ بیتی استفاده می کنیم که به صورت زیر تعریف شده است:

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

و مقادیر درهم شهرها را می توان در جدول زیر یافت.

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

همان طور که در جدول مشخص است، شهرهای لندن و مادرید مقدار درهم یکسانی دارند، که البته همان طور که ذکر شد چنین تصادم هایی مورد انتظار و کاملا طبیعی هستند. داده ساختار نمای زیر را دارد:

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

با توجه به الگوریتم شمارش خطی، کسر V بیت های خالی را در $LINEARCOUNTER$ محاسبه می کنیم:

$$V = \frac{9}{16} = 0.5625$$

و تعداد منحصربه‌فرد تخمین زده شده برابر $9.206 \approx -16 \times \ln 0.5625 \approx n$ است.

در ادامه چند ویژگی الگوریتم شمارنده خطی را برمی‌شماریم.

در صورتی که تابع درهم‌ساز h دارای زمان محاسبه‌ای ثابت باشد (که برای غالب توابع درهم‌ساز صادق است)، زمان پردازش هر مقدار برابر مرتبه $O(N)$ است، که در آن N تعداد کل مقادیر از جمله تکراری‌ها است. بنابراین، الگوریتم دارای پیچیدگی زمانی $O(N)$ است.

مانند هر الگوریتم احتمالاتی دیگر، تعدادی پارامتر وجود دارد که می‌توان برای تأثیرگذاری بر عملکرد آن تنظیم کرد.

دقت مورد انتظار تخمین به اندازه آرایه بیتی m و نسبت آن به تعداد مقادیر متمایز $\alpha = \frac{m}{n}$ که ضریب بار نامیده می‌شود، بستگی دارد. مگر اینکه $\alpha \geq 1$ ($m > n$) که در عمل جالب نیست، احتمالی غیرصفر pr اشباع وجود دارد که آرایه بیتی پر شود، که احتمال پر شدن یا اشباع نامیده می‌شود و نتایج الگوریتم را کاملاً مخدوش کرده و مقدار معادله (۱) را بی‌نهایت می‌کند. احتمال اشباع pr به ضریب بار و در نتیجه به اندازه m بستگی دارد. پس باید به اندازه کافی بزرگ انتخاب شود تا احتمال پر شدن ناچیز باشد.

خطای استاندارد δ معیاری از تغییرپذیری تخمین حاصل از شمارش خطی است که تعادلی بین آن و اندازه آرایه بیتی m وجود دارد. کاهش خطای استاندارد منجر به تخمین‌های دقیق‌تر می‌شود، اما حافظه موردنیاز را افزایش می‌دهد.

جدول ۲. تعادل 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

وابستگی به انتخاب m بسیار پیچیده است و راه‌حل تحلیلی ندارد. با این حال، نویسندگان الگوریتم مقادیر از پیش محاسبه شده‌ای برای احتمال قابل قبول تکمیل pr اشباع $= 0.7\%$ را معرفی کرده‌اند و می‌توان از آنها به عنوان مرجع استفاده کرد. به دلیل صفر نبودن احتمال اشباع، آرایه بیتی به ندرت پر می‌شود و الگوریتم ۲ را مخدوش می‌کند. هنگام کار با مجموعه‌های داده کوچک، می‌توان همه مقادیر را با تابع درهم‌ساز متفاوت دوباره نمایه کرد یا اندازه شمانده خطی را افزایش داد. اما، چنین راه‌حلی برای مجموعه‌های داده عظیم کار نخواهند کرد و همراه با پیچیدگی زمانی نسبتاً بالا، نیاز به جستجوی جانشین‌ها دارند.

با این حال، شمارش خطی زمانی که تعداد منحصربه‌فرد مجموعه داده‌ای که اندازه‌گیری می‌شود خیلی بزرگ نیست، بسیار خوب عمل می‌کند و می‌توان از آن برای بهبود الگوریتم‌های دیگر استفاده کرد که برای ارائه بهترین رفتار ممکن برای تعداد منحصربه‌فردهای عظیم ایجاد شده‌اند. در الگوریتم شمارش خطی، تخمین تعداد منحصربه‌فرد تقریباً متناسب با مقدار دقیق است، به همین دلیل از

اصطلاح «خطی» استفاده می‌شود. در بخش بعدی، الگوریتم جانسینی را در نظر می‌گیریم که می‌تواند به عنوان شمارش «لگاریتمی» طبقه‌بندی شود زیرا مبتنی بر تخمین‌هایی است که لگاریتم تعداد منحصر به فرد واقعی هستند.

شمارش احتمالاتی

فیلیپ فلاژوله و جی. نایجل مارتین در سال ۱۳۶۴ الگوریتم شمارشی بر اساس مشاهده الگوهای متداول در نمایش‌های درهم مقدار ورودی معرفی کردند. مطابق موارد قبلی، تابع درهم‌ساز h بر مقدار ورودی اعمال می‌شود و آن را به عددی از اعداد صحیح با توزیع نزدیک-یکنواخت در بازه $\{0, 1, \dots, 2^\ell - 1\}$ یا به طور معادل، به رشته‌ای از مجموعه رشته‌های دودویی به طول ℓ تبدیل می‌کند:

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

بار k

فلاژوله و مارتین تشخیص دادند که احتمال حضور الگوهای $10^k = 100\dots 0$ در رشته‌های دودویی برابر است و اگر برای هر مقدار نمایه شده ثبت شوند، می‌توانند نقش تخمین‌گر تعداد منحصر به فرد را ایفا کنند. هر الگو را می‌توان با شاخص آن، به نام «رتبه»، مرتبط کرد که از معادله زیر به دست می‌آید:

$$\text{رتبه}(i) = \begin{cases} \min_{i_k \neq 0} k, & \forall i > 0 \\ \ell, & i = 0 \end{cases} \quad (3)$$

و اساساً متناظر موقعیت سمت راست‌ترین ۱ است که به عنوان ۱ با کمترین اهمیت شناخته می‌شود.

مثال ۴- محاسبه رتبه

عدد صحیح ۸ بیتی ۴۲ را در نظر می‌گیریم که نمایش دودویی آن با طرح شماره‌گذاری ۰ LSB برابر است با:

$$42 = 0 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 = (00101010)_2$$

یک‌ها در موقعیت‌های ۱، ۳ و ۵ ظاهر می‌شوند. با توجه به تعریف در معادله (۳)، رتبه ۴۲ برابر است با

$$\text{رتبه}(42) = \min(1, 3, 5) = 1$$

رخدادهای الگوی 10^k ، یا متناظر $k = \text{رتبه}(0)$ ، در نمایش‌های دودویی مقادیر درهم‌ساز هر مقدار نمایه شده را می‌توان به صورت فشرده در داده‌ساختار آرایه‌ای شمارنده به طول ℓ ذخیره شوند و یا «طرح فم» نیز شناخته می‌شود.

در ابتدا، همه بیت‌ها در شمارنده برابر با صفر هستند. هنگام درج مقدار جدید x در داده‌ساختار، مقدار درهم آن را با تابع درهم‌ساز h محاسبه، سپس رتبه متناظر را محاسبه، و بیت متناظر را در آرایه برابر یک قرار می‌دهیم. مراحل در الگوریتم زیر بیان شده است.

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 $COUNTER[j] = 0$

$COUNTER[j] \leftarrow 1$

به این ترتیب، یک در شمارنده در موقعیت j به معنای حداقل یک بار مشاهده الگوی 10^j در میان مقادیر درهم شده همه مقادیر نمایه شده است.

```

الگوریتم ۴- محاسبه رتبه X
ورودی: عدد صحیح X
خروجی: رتبه X یا کوچکترین j با شرط بیت j-ام برابر ۱
    j ← ۰
    l = len(x)
    تا زمانی که j < l که j ادامه بده
    اگر (1 & j >> x) == ۱
        برگرداندن j
    j ← j + ۱
برگرداندن l // ۱ رخ ندادن به دلیل غیر صفر بودن h(x)
    
```

مثال ۵- ساخت شمارنده: مجموعه داده بیست نام از پایتخت‌ها اعم از برلین، برلین، برلین، پاریس، برلین، لیسبون، کیف، پاریس، لندن، رم، آتن، مادرید، وین، رم، رم، لیسبون، برلین، پاریس، لندن، کیف، واشنگتن را در نظر می‌گیریم. از تابع درهم مورمور ۳ نوع ۳۲ بیتی به عنوان تابع درهم‌ساز h استفاده می‌کنیم که مقادیر را به $\{0, 1, \dots, 2^{32} - 1\}$ نگاشت می‌کند، در نتیجه از شمارنده‌ای به طول $l = 32$ بهره می‌گیریم. با استفاده از مقادیر درهم‌ساز محاسبه شده در مثال ۳ و تعریف (۳)، رتبه‌ها را برای هر مقدار محاسبه می‌کنیم:

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

بنابراین، شمارنده شکل زیر را دارد:

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 فرض توزیع یکنواخت مقادیر در ادامه بر اساس مشاهدات نتیجه‌ای نظری را مطرح می‌کنیم. در صورتی که تعداد مقادیر متمایز تاکنون نمایه شده برابر n باشد، انتظار می‌رود که مقدار ۱ در موقعیت اول در حدود $\frac{n}{p}$ دفعه ظاهر شود، ظاهر شدن مقدار یک در موقعیت دوم در حدود $\frac{n}{p^2}$ دفعات، و به همین ترتیب تعداد رخدادها در هر موقعیت به دست م‌آید. تعداد ظهور یک در موقعیت k برابر $\frac{n}{p^k}$ است. اگر $\log_p n \gg j$ باشد، احتمال مشاهده مقدار ۱ در موقعیت j -ام به صفر میل می‌کند. در نتیجه مقدار $\text{Counter}[j]$

تقریباً با قطعیت صفر خواهد بود. به طور مشابه، اگر $\log_2 n \ll j$ باشد، مقدار Counter[j] تقریباً با قطعیت برابر با یک خواهد بود. اگر مقدار j در حدود $\log_2 n$ باشد، احتمال مشاهده ۱ یا ۰ در آن موقعیت تقریباً برابر است.

در نتیجه، پس از درج تمام مقادیر مجموعه داده در آرایه شمارنده، می توان سمت راست ترین موقعیتی (R) را که مقدار صفر ظاهر می شود به عنوان شاخص تخمین $\log_2 n$ در نظر گرفت. فلاژوله و مارتین نشان دادند که در عمل، ضریب تصحیح ϕ جهت بهبود تخمین لازم است و تخمین کاردینالیته (تعداد مقادیر متمایز) با معادله زیر انجام می شود:

$$n \approx \frac{1}{\phi} 2^R \quad (4)$$

که $\phi \approx 0.77351$ است. سخن کوتاه، فلاژوله و مارتین از موقعیت کم ارزش ترین بیت صفر (راست ترین صفر) به عنوان تخمین کاردینالیته استفاده کردند و الگوریتم خود را بر همین اساس ساختند. با این حال، از مشاهده بالا می توان دید که موقعیت پر ارزش ترین بیت ۱ (چپ ترین ۱) نیز می تواند به همین منظور استفاده شود. اما توزیع آن یکنواخت تر است و در نتیجه خطای استاندارد بزرگ تری ایجاد می کند. الگوریتم محاسبه راست ترین موقعیت صفر در شمارنده ساده را می توان به صورت زیر معادله بندی کرد.

الگوریتم ۵- محاسبه موقعیت سمت چپ ترین صفر
ورودی: شمارنده ساده به طول l
خروجی: سمت چپ ترین موقعیت صفر
برای j از ۰ تا $l - 1$ انجام بده
اگر $COUNTER[j] = 0$ آن گاه
برگرداندن j
برگرداندن l

مثال ۶: تخمین تعداد منحصر به فرد با شمارنده ساده

آرایه مثال ۵ را در نظر می گیریم و تعداد تخمینی مقادیر منحصر به فرد را محاسبه کنید.

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 ۵، در COUNTER چپ ترین مقدار ۰ در موقعیت $R = 4$ ظاهر می شود، بنابراین، با توجه به معادله (۴)، تخمین تعداد منحصر به فرد است

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

تعداد دقیق منحصر به فرد مجموعه ۱۰ است، در نتیجه حاصل تخمین دارای خطای زیادی است زیرا مقادیر R اعداد صحیح هستند و برای رتبه های بسیار نزدیک می توانیم نتایجی را به دست آوریم که دارای تفاوتی از مرتبه توان های دودویی باشند. در مثال حاضر نیز، $R=3$ تخمین نسبتاً مناسب ۱۰.۳۴ را تولید می کند.

تخمین تعداد منحصر به فرد بر اساس شمارنده معمول می‌تواند مقادیر مورد انتظار بسیار نزدیکی را ارائه دهد، اما دارای بردایی نسبتاً بالایی است که معمولاً، همانطور که در مثال قبل نیز مشاهده کردیم، با خطای استاندارد δ دارای اختلاف از مرتبه دو برابر است.

ناگفته پیداست، ضعف رویکرد تک‌شمارنده فقدان تخمین‌های با اطمینان بالا برای تعداد منحصر به فرد است. در واقع، پیش‌بینی خود را تنها بر اساس تخمینی واحد انجام می‌دهد. از این رو، بهبود منطقی الگوریتم استفاده از تعداد زیادی شمارنده و در نتیجه افزایش تعداد تخمین‌ها است. پیش‌بینی نهایی n را می‌توان با میانگین‌گیری از پیش‌بینی‌های R_i از شمارنده‌های مزبور $\left\{ \text{شمارنده}_i \right\}_{i=1}^{m-1}$ به دست آورد. سخن کوتاه، معادله اصلاح شده (۴) الگوریتم شمارش احتمالاتی به شکل زیر است:

$$n \approx \frac{1}{\phi} \frac{1}{m} \sum_{i=1}^{m-1} R_i \quad (5)$$

و تعداد منحصر به فرد n مقدار تخمینی با همان کیفیت اما با بردایی بسیار کوچکتر را در پی خواهد داشت. اشکال عملی برای ساخت m شمارنده مستقل، نیاز به محاسبه مقادیر m تابع درهم‌ساز متفاوت است که با توجه به اینکه هر تابع درهم‌ساز را می‌توان در $O(1)$ محاسبه کرد، محاسبه همه آنها دارای پیچیدگی زمانی $O(m)$ و در نتیجه هزینه‌های پردازشی نسبتاً بالا را در پی دارد. راه‌حل بهینه‌سازی الگوریتم شمارش احتمالاتی، اعمال «میانگین‌گیری تصادفی» است. به دیگر سخن، m تابع درهم‌ساز صرفاً با یک تابع جانشین می‌شوند اما حاصل مقدار آن به دو بخش خارج‌قسمت و باقیمانده تقسیم می‌شود. باقیمانده r برای انتخاب یکی از m شمارنده و خارج‌قسمت q برای محاسبه رتبه و یافتن شاخص مناسب برای به‌روزرسانی در آن شمارنده استفاده می‌شود.

الگوریتم ۶: استفاده از میانگین‌گیری برای به‌روزرسانی شمارنده‌ها

ورودی: مقدار $x \in D$

ورودی: آرایه‌ای از m شمارنده ساده با تابع درهم‌ساز h

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

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

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

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

با اعمال الگوریتم میانگین‌گیری تصادفی ۶ به شمارش احتمالاتی، با این فرض که توزیع مقادیر مبتنی بر خارج‌قسمت به اندازه کافی منصفانه است، ممکن است انتظار داشته باشیم که $\frac{n}{m}$ مقدار توسط هر شمارنده ساده $\{\text{COUNTER}_i\}_{i=1}^{m-1}$ اندیس شده باشند، بنابراین معادله (۵) تخمین خوبی برای $\frac{n}{m}$ و نه مستقیماً n است. پس معادله‌ای که فلاژوله و مارتین پیشنهاد دادند به صورت زیر و الگوریتم پیشنهادی آنها به صورت زیر است:

$$n \approx \frac{m}{\phi} \bar{r} = \frac{m}{\phi} \frac{1}{m} \sum_{i=1}^{m-1} R_i \quad (6)$$

الگوریتم ۷: الگوریتم فلاژوله-مارتین (PCSA)

ورودی: مجموعه داده D

ورودی: آرایه‌ای از m شمارنده ساده با تابع درهم‌ساز h

خروجی: تخمین تعداد منحصر به فرد

برای $x \in D$ انجام بده

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

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

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

$$COUNTER_{f_r}[j] \leftarrow 1$$

برای r از ۰ تا $m - 1$ انجام بده

$$R \leftarrow \text{LeftMostZero}(COUNTER_{f_r})$$

$$S \leftarrow S + R$$

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

برگرداندن $\frac{m}{\phi} 2^S$

الگوریتم بالا الگوریتم شمارش احتمالاتی با میانگین‌گیری تصادفی (PCSA) نامیده می‌شود و همچنین الگوریتم فلاژوله-مارتین نیز شناخته می‌شود. در مقایسه با نسخه‌ای که از m تابع درهم‌ساز استفاده می‌کند، پیچیدگی زمانی را برای هر مقدار به حدود $O(1)$ کاهش می‌دهد.

مثال ۷- تخمین تعداد منحصربه‌فرد با میانگین‌گیری تصادفی: مجموعه داده و مقادیر درهم‌ساز محاسبه شده در مثال ۶ را در نظر می‌گیریم و فن میانگین‌گیری تصادفی را با شبیه‌سازی $m = 3$ تابع درهم‌ساز اعمال می‌کنیم. از باقیمانده r برای انتخاب یکی از سه شمارنده و از خارج‌قسمت q برای محاسبه رتبه استفاده می‌کنیم.

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

هر شمارنده اطلاعات مربوط به حدود یک سوم شهرها را مدیریت می‌کند، بنابراین، توزیع به اندازه کافی منصفانه است. پس از نمایه کردن همه مقادیر و تنظیم بیت‌های مناسب در شمارنده‌های مربوطه، شمارنده‌های حاصل به دست می‌آیند. شمارنده ۰ شکل زیر را دارد:

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

شمارنده ۱ به صورت زیر است:

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

شمارنده ۲ به صورت زیر است:

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

سمت چپ‌ترین موقعیت‌های صفر برای هر شمارنده (که در بالا مشخص شده‌اند) عبارت از $R_1 = 3$ ، $R_2 = 1$ و $R_3 = 1$ هستند. بنابراین، تخمین تعداد منحصربه‌فرد با توجه به معادله (۶) است

$$n \approx \frac{3}{\phi} \frac{1}{2^{\sum_{k=1}^3 R_k}} \approx \frac{3}{0.77351} 2^{-\frac{1+2+1}{3}} \approx 12.31$$

تخمین حاصل بسیار نزدیک‌تر به مقدار واقعی تعداد منحصربه‌فرد ۱۰ شده است، و حتی بدون استفاده از شمارنده‌های زیاد، به طور قابل توجهی از تخمین مثال ۶ بهتر عمل می‌کند.

ویژگی‌ها

الگوریتم فلاژوله-مارتین برای مجموعه‌های داده با تعداد منحصربه‌فرد زیاد به خوبی کار می‌کند و در نسبت‌های $\frac{n}{m} > 20$ تقریب‌های خوبی به دست می‌دهد. با این حال، می‌توان از الگوریتم‌های غیرخطی در الگوریتم برای کار با مجموعه داده‌های با تعداد منحصربه‌فردهای کوچک بهره برد تا نتایج را بهبود داد. بیورن شویرمان و مارتین ماوه در سال ۱۳۸۶ اصلاحی احتمالی پیشنهاد کردند که و معادله (۶) را با افزودن عبارتی برای تعداد منحصربه‌فردهای کوچک متناسب ساختند که برای تعداد منحصربه‌فردهای بزرگ به صفر همگرا می‌شود:

$$n \approx \frac{m}{\phi} (2^{\bar{R}} - 2^{-\lambda \bar{R}}) \quad (7)$$

که $\lambda \approx 1.75$ خطای استاندارد δ الگوریتم فلاژوله-مارتین با تعداد شمارنده‌های استفاده شده رابطه معکوس دارد و می‌توان آن را به صورت تقریبی محاسبه کرد:

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

مقادیر مرجع خطای استاندارد برای تعداد شمارنده‌های پرکاربرد را می‌توان در جدول ۵ یافت. طول ℓ هر شمارنده را می‌توان به گونه‌ای انتخاب کرد که:

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

بنابراین، $\ell = 32$ عملاً برای شمارش تعداد منحصربه‌فردهای بسیار فراتر از 10^9 با استفاده از ۶۴ شمارنده کافی است.

جدول ۵- تعادل بین دقت و ذخیره‌سازی ($\ell = 32$)

δ	حافظه	m
۹.۷ درصد	۲۵۶ بایت	۶۴
۴.۸ درصد	۱۰۲۴ کیلوبایت	۲۵۶
۲.۴ درصد	۴.۱ کیلوبایت	۱۰۲۴

شمارنده‌های متناظر مجموعه‌های داده مختلف را می‌توان به راحتی با اعمال عملیات OR بیتی با هم ادغام کرد که نتیجه آن یک شمارنده برای کل آن مجموعه داده‌هاست. مانند فیلتر بلوم، الگوریتم‌های شمارش احتمالاتی از حذف پشتیبانی نمی‌کنند. اما، از رویکردی مشابه بلوم شمارنده و با افزایش فضا می‌توان حذف را پوشش داد.

روش‌های لگ‌لگ و ابرلگ‌لگ

پرکاربردترین الگوریتم‌های احتمالاتی برای تخمین تعداد منحصر به فرد که در عمل استفاده می‌شوند، خانواده لگ‌لگ از الگوریتم‌ها است که شامل الگوریتم لگ‌لگ LogLog است که پیشنهادی ماریان دوراند و فیلیپ فلاژوله در سال ۱۳۸۲، و دو روش جانشین آن ابرلگ‌لگ و ابرلگ‌لگ++ است. الگوریتم‌ها از رویکردی مشابه الگوریتم شمارش احتمالاتی استفاده می‌کنند. به بیان دیگر، تخمین تعداد منحصر به فرد n با مشاهده حداکثر تعداد صفرهای ابتدایی در نمایش دودویی مقادیر انجام می‌شود. همه آنها به حافظه کمکی نیاز دارند و داده‌ها را یکبار پیمایش می‌کنند تا تخمینی از تعداد منحصر به فرد تولید کنند.

مطابق منوال روش‌های قبلی تعیین تعداد منحصر به فرد اجزا، هر مقدار در مجموعه داده با اعمال تابع درهم‌سازی پیش‌پردازش می‌شود که مقادیر را 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 M تبدیل می‌کند:

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

مراحل الگوریتم‌ها مشابه PCSA است که در اینجا بار دیگر یادآوری می‌کنیم. ابتدا، مجموعه داده اولیه یا جریان ورودی را به تعدادی زیرمجموعه تقسیم می‌کند که هر کدام با یکی از m شمارنده ساده نمایه می‌شوند. سپس، با توجه به میانگین‌گیری تصادفی، زیرا تابع درهم‌ساز واحد وجود دارد، شمارنده را برای مقدار خاص x با استفاده از یک بخش از مقدار درهم $h(x)$ انتخاب می‌کنیم، در حالی که بخش دیگر برای به‌روزرسانی شمارنده مربوطه استفاده می‌شود.

همه الگوریتم‌های مورد بحث در اینجا بر اساس مشاهده الگوهای 10^k هستند که در ابتدای مقادیر برای شمارنده خاص رخ می‌دهند، و هر الگو را با شاخص آن، به نام رتبه، مرتبط می‌کنند. رتبه معادل موقعیت ۱ با کمترین اهمیت در نمایش دودویی مقدار درهم‌ساز مقدار نمایه شده است و می‌تواند با معادله (۳) محاسبه شود. هر شمارنده مشاهده تعداد منحصر به فرد خود را بر اساس رتبه‌های دیده شده می‌سازد، و تخمین نهایی تعداد منحصر به فرد از چنین مشاهداتی با استفاده از تابع ارزیابی تولید می‌شود.

در مورد ذخیره‌سازی، نگهداری شمارنده‌ها در الگوریتم شمارش احتمالاتی نسبتاً پرهزینه است، اما الگوریتم لگ‌لگ راه‌حلی فضا-کارآمدتر را همراه با تابع ارزیابی بهتر و رویکرد تصحیح سوگیری پیشنهاد می‌کند.

الگوریتم لگ‌لگ

الگوریتم لگ‌لگ بر اساس محاسبه رتبه هر مقدار ورودی با تابع درهم‌ساز واحد h قرار دارد. در صورتی که n تعداد کل عنصرهای نمایه شده در شمارنده است، می‌توان ادعا کرد که $\frac{n}{2^k}$ مقدار رتبه k اختیار کنند پس رتبه بیشینه مشاهده شده نشانه و تقریب مناسبی از مقدار $\log_2 n$ فراهم کند:

$$R = \max_{x \in D} (\text{rank}(x)) \approx \log_2 n \quad (8)$$

ولی تخمین مذکور دارای خطایی با مرتبه ۱.۸۷ برابر مقدار واقعی بود که عملاً بی‌معنی است. برای کاهش خطا، الگوریتم لگ‌لگ از سطل‌بندی مبتنی بر میانگین‌گیری تصادفی استفاده می‌کند و مجموعه داده را به $m = 2^p$ زیرمجموعه $\{S_0, S_1, \dots, S_{m-1}\}$ تقسیم می‌کند، که در آن پارامتر دقت p تعداد بیت‌های استفاده شده را تعریف می‌کند. بنابراین، برای هر مقدار x از مجموعه داده، p بیت اول مقدار درهم h -بیتی $h(x)$ را می‌توان برای یافتن شاخص j زیرمجموعه مناسب گرفت:

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

و بقیه $(\ell - p)$ بیت در شمارنده متناظر $COUNTER[j]$ اندیس می‌شوند تا رتبه محاسبه شده و مشاهده R_j طبق معادله (۸) به دست آید.

با توزیع منصفانه، هر زیرمجموعه دارای n/m عضو خواهد بود. پس، R_j -ها در شمارنده‌ها $\{COUNTER[j]\}_{j=0}^{m-1}$ می‌توانند نشانه‌ای از مقدار $\log_2 \frac{n}{m}$ ارائه دهند، و با استفاده از میانگین حسابی آنها و تصحیح سوگیری، می‌توان وردائی مشاهده را کاهش داد:

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

که در آن $\alpha_m = \left(\Gamma\left(-\frac{1}{m}\right) \times \frac{1 - 2^{\frac{1}{m}}}{\log 2} \right)^m \approx 0.39701 - \frac{2\pi^2 + (ln 2)^2}{48m}$ و $\Gamma(\cdot)$ تابع گاما است. با این حال، برای اکثر موارد عملی $m \geq 64$ استفاده از $\alpha_m \approx 0.39701$ کفایت می‌کند.

الگوریتم ۸- تخمین تعداد منحصره‌فرد با لگ‌لگ
 ورودی: مجموعه داده D
 ورودی: آرایه‌ای از m شمارنده لگ‌لگ با تابع درهم‌ساز h
 خروجی: تخمین تعداد منحصره‌فرد

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

برای $x \in D$ انجام بده

$i \leftarrow h(x) = (i_{\ell-1} \dots i_1 i_0)_p, i_k \in \{0, 1\}$
 $j \leftarrow (i_{p-1} \dots i_1 i_0)_p$
 $r \leftarrow rank\left(\left(i_{\ell-1} \dots i_{p+1} i_p\right)_p\right)$
 $COUNTER[j] \leftarrow \max(COUNTER[j], r)$
 $R \leftarrow \frac{1}{m} \sum_{j=0}^{m-1} COUNTER[j]$
 برگرداندن $\alpha_m \cdot m \cdot 2^R$

مثال ۸- فرض کنید $m=4$ پس در نتیجه $p = \log_2 4 = 2$ بیت جهت انتخاب شمارنده نیاز داریم. فرض می‌کنیم خروجی درهم طولی برابر هشت بیت دارد. پس بیت‌های صفر و یک برای شمارنده و بقیه شش بیت برای محاسبه رتبه استفاده می‌شود.

شهر	h(x) (^ bits)	(b ₁ b ₀)	(b ₇ b ₆ b ₅ b ₄ b ₃ b ₂)
Tehran	01101010	10 = 2	011010
Paris	10110011	11 = 3	101100
London	00011100	00 = 0	000111
Berlin	11001010	10 = 2	110010

001011	01 = 1	00101101	Rome
010011	10 = 2	01001110	Madrid
100001	11 = 3	10000111	Kyiv
111000	00 = 0	11100000	Vienna

در قدم بعدی رتبه بیت‌های بالایی حساب می‌شود.

رتبه	$(b_7b_6b_5b_4b_3b_2)$	شهر
1	011010	Tehran
2	101100	Paris
0	000111	London
1	110010	Berlin
0	001011	Rome
0	010011	Madrid
0	100001	Kyiv
3	111000	Vienna

در قدم بعدی شمارنده‌ها با عملیات بیش بروز می‌شوند.

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

پس شمارنده دارای مقدار $[3,0,1,2]$ است. میانگین را به صورت زیر حساب می‌کنیم.

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

نتیجه نهایی تخمین برای $m=4$ و $\alpha_4 \approx 0.67$ برابر است با

$$\hat{n} = \alpha_m \times m \times 2^R = 0.67 \times 4 \times 2^{1.5} = 7.58$$

تعداد واقعی برابر با 8 است و تخمین لگ‌لگ برابر 7.6 است و خطا در حدود پنج درصد است.

ویژگی‌ها

خطای استاندارد δ الگوریتم لگ‌لگ با تعداد شمارنده‌های استفاده شده m رابطه معکوس دارد و می‌توان آن را به صورت تقریبی محاسبه کرد:

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

بنابراین، برای $m = 256$ خطای استاندارد حدود هشت درصد و برای $m = 1024$ به حدود چهار درصد کاهش می‌یابد.

اگر شمارش تا n لازم باشد، فضای مورد نیاز ذخیره‌سازی الگوریتم لگ‌لگ تخمینی از $O(\log_2 \log_2 n)$ است. به طور دقیق‌تر، کل

فضای مورد نیاز الگوریتم برای شمارش تا n برابر $m \cdot \log_2 \log_2 \frac{n}{m} (1 + O(1))$ است.

در مقایسه با الگوریتم شمارش احتمالاتی که در آن هر شمارنده به ۱۶ یا ۳۲ بیت نیاز دارد، الگوریتم لگالگ به شمارنده‌های بسیار کوچکتر $\{COUNTER[j]\}_{j=0}^{m-1}$ معمولاً هر کدام در حدود پنج بیت نیاز دارد. با وجود مصرف حافظه کمتر الگوریتم لگ نسبت به الگوریتم شمارش احتمالاتی، به میزان ناچیزی دقت کاهش می‌یابد.

در صورتی که به‌شمارش تعداد منحصره‌فرد تا حدود 2^{30} ، یعنی حدود ۱ میلیارد، با دقت حدود چهار درصد نیاز باشد. همانطور که گفته شد، برای چنین خطای استاندارد، $m = 1024$ سطل مورد نیاز است که هر کدام حدود $\frac{n}{m} = 2^{20}$ مقدار دریافت می‌کنند. 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) کافی است. بنابراین، برای تخمین تعداد منحصره‌فرد تا حدود 10^9 با خطای استاندارد چهار درصد، الگوریتم به 1024 سطل ۵ بیتی نیاز دارد که در مجموع ۶۴۰ بیت است.

الگوریتم ابرلگالگ

فیلیپ فلاژوله، اریک فیوزی، اولیویه گاندوئه، و فردریک مونیه در سال ۱۳۸۶ الگوریتم ابرلگالگ را بر اساس الگوریتم لگالگ پیشنهاد دادند. الگوریتم ابرلگالگ از تابع درهم‌ساز ۳۲ بیتی، تابع ارزیابی متفاوت و تصحیحات سوگیری مختلف استفاده می‌کند. مشابه الگوریتم لگالگ، ابرلگالگ از تصادفی‌سازی برای تقریب تعداد منحصره‌فرد مجموعه داده استفاده می‌کند و برای مدیریت تعداد منحصره‌فرد تا 10^9 با یک تابع درهم‌ساز ۳۲ بیتی h طراحی شده است و مجموعه داده را به $m = 2^p$ زیرمجموعه، با دقت $p \in \{16, \dots, 4\}$ تقسیم می‌کند.

علاوه بر این، تابع ارزیابی الگوریتم ابرلگالگ را از لگالگ استاندارد متمایز می‌کند. الگوریتم لگالگ اصلی از میانگین هندسی استفاده می‌کند در حالی که ابرلگالگ از تابعی استفاده می‌کند که مبتنی بر نسخه نرمال شده میانگین هارمونیک است:

$$\hat{n} \approx \alpha_m \cdot m^2 \cdot \left(\sum_{j=0}^{m-1} \gamma^{-COUNTER[j]} \right)^{-1} \quad (9)$$

که در آن

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

مقادیر تقریبی α_m را در جدول زیر یافت. شهود استفاده از میانگین هارمونیک این است که وردائی را به دلیل ویژگی آن در هموارسازی توزیع‌های احتمال چوله کاهش می‌دهد.

جدول ۶- α_m برای پرکاربردترین مقادیر m

m	α_m
2^4	۰.۶۷۳
2^5	۰.۶۹۷
2^6	۰.۷۰۹
$\geq 2^7$	$\frac{0.7213 \times m}{m + 1.079}$

تخمین (۹) به دلیل خطاهای غیر خطی نیاز به تصحیح برای محدوده‌های کوچک و بزرگ دارد. فلاژوله و همکاران در نتیجه آزمایش تجربی دریافتند که برای تعداد منحصربه‌فردهای کوچک $n < \frac{5}{3}m$ برای دستیابی به تخمین‌های بهتر، الگوریتم ابرلگ را می‌توان با شمارش خطی با استفاده از تعدادی از شمارنده‌های COUNTER[j] غیر صفر تصحیح کرد (اگر یک شمارنده مقدار صفر داشته باشد، می‌توانیم با قطعیت بگوییم که آن زیرمجموعه خاص خالی است). بنابراین، برای محدوده‌های مختلف تعداد منحصربه‌فرد، که به صورت فواصل روی تخمین \hat{n} محاسبه شده توسط معادله (۳.۱۲) بیان می‌شوند، الگوریتم تصحیحات زیر را ارائه می‌دهد:

$$n = \begin{cases} \text{LinearCounter}, \hat{n} \leq \frac{5}{3}m \wedge \exists j: \text{Counter}[j] \neq 0 \\ -2^{23} \log\left(1 - \frac{\hat{n}}{2^{23}}\right), & \hat{n} > \frac{1}{3} \cdot 2^{23} \\ \hat{n}, & \text{دغا} \end{cases}$$

با این حال، برای $n = 0$ تصحیح به درستی کار نمی‌کند و الگوریتم همیشه تقریباً $0.7m$ را برمی‌گرداند.

به دلیل استفاده الگوریتم ابرلگ از تابع درهم‌ساز ۳۲ بیتی، هنگامی که تعداد منحصربه‌فرد به $4 \times 10^9 \approx 2^{32}$ نزدیک می‌شود، تابع درهم‌ساز تقریباً به حد خود می‌رسد و احتمال تصادم افزایش می‌یابد. برای چنین محدوده‌های بزرگ، الگوریتم ابرلگ تعداد مقادیر درهم‌ساز مختلف را تخمین می‌زند و از آن برای تقریب تعداد منحصربه‌فرد استفاده می‌کند. با این حال، در عمل، این خطر وجود دارد که تعداد بالاتر به سادگی قابل نمایش نباشد و از بین برود و بر دقت تأثیر بگذارد.

تابع درهم‌ساز را در نظر می‌گیریم که دامنه را به مقادیر M بیت نگاشت می‌کند. چنین تابعی حداکثر می‌تواند 2^M مقدار مختلف را رمزگذاری کند و اگر تعداد منحصربه‌فرد تخمینی n به چنین حدی نزدیک شود، تصادم‌های درهم‌ساز محتمل‌تر می‌شوند.

هیچ نشانی دال بر عملکرد بهتر توابع درهم‌ساز پرکاربرد (مانند MurmurHash^3 , MD^5 , SHA-1 , SHA-256) نسبت به سایر روش‌ها در الگوریتم‌های ابرلگ یا اصلاحات آن وجود ندارد.

الگوریتم کامل ابرلگ در زیر نشان داده شده است.

```

الگوریتم ۹- تخمین تعداد منحصربه‌فرد با ابرلگ
ورودی: مجموعه داده D
ورودی: آرایه‌ای از m شمارنده LogLog با تابع درهم‌ساز h
خروجی: تخمین تعداد منحصربه‌فرد
COUNTER[j] ← 0, j = 0 ... m - 1
برای x ∈ D انجام بده
    i ← h(x) = (i_{31} ... i_1 i_0)_2
    j ← (i_{p-1} ... i_1 i_0)_2
    r ← rank((i_{31} ... i_{p+1} i_p)_2)
COUNTER[j] ← max(COUNTER[j], r)
R ← ∑_{j=0}^{m-1} 2^{-Counter[j]}
n̂ = α_m · m^2 · 1/R
n ← n̂
    
```

اگر $\hat{n} \leq \frac{5}{4}m$ آنگاه
 $Z \leftarrow \text{تعداد } j=0, \dots, m-1 \text{ (COUNTER}[j] = 0)$
 اگر $Z \neq 0$ آنگاه
 $n \leftarrow m \cdot \log\left(\frac{m}{Z}\right)$
 در غیر این صورت اگر $\hat{n} > \frac{1}{3} \cdot 2^{32}$ آنگاه
 $n \leftarrow -2^{32} \cdot \log\left(1 - \frac{n}{2^{32}}\right)$
 برگرداندن n

مثال ۹- در مثال قبل مقدار شمارنده‌ها دارای مقدار $[2, 1, 0, 3]$ است. میانگین را به صورت زیر حساب می‌کنیم.

$$R = \sum_{j \text{ شمارنده}} 2^{-j} = (2^{-3} + 2^0 + 2^{-1} + 2^{-2}) = \frac{1}{8} + 1 + \frac{1}{2} + \frac{1}{4} = \frac{15}{8} = 1.875$$
 نتیجه نهایی تخمین برای $m=4$ و $\alpha_4 \approx 0.673$ برابر است با

$$\hat{n} = \alpha_m \times m^2 \times \frac{1}{R} = 0.673 \times 16 \times \frac{1}{1.875} = 5.74$$
 تعداد واقعی برابر با ۸ است و تخمین لگ‌لگ برابر ۵.۷۴ است.

ویژگی‌ها

مشابه الگوریتم لگ‌لگ، بین خطای استاندارد δ و تعداد شمارنده‌ها m نیاز به سبک-سنگینی وجود دارد:

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

نیاز حافظه با تعداد مقادیر به صورت خطی رشد نمی‌کند (برخلاف الگوریتم‌هایی مانند الگوریتم شمارش خطی)، با تخصیص $(M-p)$ بیت برای مقادیر درهم‌ساز و داشتن $m = 2^p$ شمارنده در مجموع، حافظه مورد نیاز برابر است با

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

به بیان دیگر، هر شمارنده به $(1 + \text{رتبه حداکثر}) \cdot \log_2$ بیت نیاز دارد. علاوه بر این، الگوریتم فقط از توابع درهم‌ساز ۳۲ بیتی و دقت ۱۶... ۴ استفاده می‌کند، نیازهای حافظه برای داده‌ساختار ابرلگ‌لگ برابر با $2^p \times 5$ بیت است.

بنابراین، الگوریتم ابرلگ‌لگ تخمین تعداد منحصربه‌فردهای بسیار فراتر از 10^9 را با دقت دو درصد و با استفاده از تنها ۱.۵ کیلوبایت حافظه ممکن می‌سازد.

به عنوان مثال، پایگاه داده معروف درون حافظه [Redis](#) ساختارهای داده ابرلگ‌لگ با ۱۲ کیلوبایت را حفظ می‌کند که تعداد منحصربه‌فردها را با خطای استاندارد ۰.۸۱ درصد تقریب می‌زند. جدول زیر مقایسه ای از ابرلگ‌لگ را با دیگر الگوریتم‌ها نشان می‌دهد.

روش	حافظه مورد نیاز برای یک میلیارد مقدار یکتا	خطا
شمارش دقیق	حدود ۳۲ گیگابایت	۰ درصد
شمارش خطی	حدود یک گیگابایت	۱ درصد
ابرلگ‌لگ	حدود ۱۲ کیلوبایت	۰.۸۱ درصد

در حالی که ابرلگ در مقایسه با لگ تخمین تعداد منحصره‌فرد را برای مجموعه‌های داده کوچک بهبود بخشید، هنوز هم تعداد منحصره‌فردهای واقعی را در چنین مواردی بیش از حد تخمین می‌زند. انواع الگوریتم‌های ابرلگ در پایگاه‌های داده شناخته شده‌ای مانند Amazon Redshift, Redis, Apache CouchDB, Riak و جز اینها پیاده‌سازی شده‌اند.

الگوریتم ابرلگ++

اشتفان هویله، مارک نونکسر، و الکساندر هال در سال ۱۳۹۲، نسخه بهبود یافته ابرلگ را به نام ابرلگ++ معرفی کردند که بر روی تعداد منحصره‌فردهای بزرگ و تصحیح سوگیری بهتر متمرکز است. اهم اصلاحات در الگوریتم ابرلگ++ استفاده از تابع درهم‌ساز ۶۴ بیتی است. واضح است که هرچه مقادیر خروجی تابع درهم‌ساز طولانی‌تر باشد، مقادیر مختلف بیشتری را می‌توان کدگذاری کرد. چنین بهبودی امکان تخمین تعداد منحصره‌فردهای بسیار بزرگتر از 10^9 مقدار منحصره‌فرد را فراهم می‌کند، اما زمانی که تعداد منحصره‌فرد به $10^{19} \cdot 1.8 \approx 2^{64}$ نزدیک می‌شود، تصادم‌های درهم برای ابرلگ++ نیز مشکل‌ساز می‌شوند.

الگوریتم ابرلگ++ عیناً از تابع ارزیابی (۹) استفاده می‌کند. با این حال، تصحیح سوگیری را بهبود می‌بخشد. نویسندگان الگوریتم مجموعه‌ای از آزمایش‌ها را برای اندازه‌گیری سوگیری انجام دادند و دریافتند که برای $n \leq 5m$ سوگیری الگوریتم ابرلگ صلی را می‌توان با استفاده از داده‌های تجربی جمع‌آوری شده در طول آزمایش‌ها، تصحیح بیشتر کرد.

علاوه بر مقاله اصلی، هویله و همکاران مقادیر تعیین شده تجربی را در آرایه‌هایی از تخمین‌های خام تعداد منحصره‌فرد داده تخمینی خام RAWESTIMATEDATA و سوگیری‌های مرتبط را در داده‌سوگیری BIASDATA برای بهبود تصحیح سوگیری در الگوریتم عرضه کردند. البته، پوشش هر مورد بالقوه ممکن نیست، در نتیجه RAWESTIMATEDATA آرایه‌ای از ۲۰۰ نقطه درونیابی را در بر دارد که میانگین تخمین خام اندازه‌گیری شده در این نقطه را در بیش از ۵۰۰۰ مجموعه داده مختلف ذخیره دارد. BIASDATA حدود ۲۰۰ سوگیری اندازه‌گیری شده را شامل می‌شود که با RAWESTIMATEDATA مطابقت دارند. هر دو آرایه صفر-اندیس هستند و حاوی مقادیر از پیش محاسبه شده برای همه دقت‌های پشتیبانی شده $p \in 4 \dots 18$ هستند، جایی که اندیس صفر در آرایه‌ها با مقدار دقت ۴ مطابقت دارد. به عنوان مثال، برای $m = 2^{10}$ و $p = 10$ داده‌های مورد نیاز را می‌توان در $BIASDATA[p]$ و $RAWESTIMATEDATA[p]$ یافت. روش تصحیح سوگیری در الگوریتم ابرلگ++ را می‌توان به صورت زیر انجام داد.

الگوریتم ۱۰- تصحیح سوگیری با استفاده از مقادیر تجربی

ورودی: تخمین \hat{n} با دقت p

خروجی: تخمین تعداد منحصره‌فرد تصحیح شده با سوگیری

$n_{low} \leftarrow 0, n_{up} \leftarrow 0, j_{low} \leftarrow 0, j_{up} \leftarrow 0$

برای $j \leftarrow 0$ تا طول($RAWESTIMATEDATA[p - 4]$) انجام بده

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

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

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

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

break

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

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

$$y = \text{درون‌یابی} \left((n_{low}, n_{low} - b_{low}), (n_{up}, n_{up} - b_{up}) \right)$$

$$y(\hat{n}) \text{ برگرداندن}$$

مثال ۱۰- تصحیح سوگیری با استفاده از مقادیر تجربی
فرض کنید که تخمین تعداد منحصربه‌فرد $\hat{n} = 2018.34$ را با استفاده از معادله (۳.۱۲) محاسبه کرده‌ایم و می‌خواهیم آن را برای دقت $(m = 2^{10})$ $p = 10$ تصحیح کنیم.
ابتدا، آرایه $RAWESTIMATEDATA[6]$ را بررسی می‌کنیم و تعیین می‌کنیم که مقدار مذکور \hat{n} بین مقادیر با اندیس‌های ۷۳ و ۷۴ آن آرایه قرار می‌گیرد، جایی که $RAWESTIMATEDATA[6][73] = 2003.1804$ و $RAWESTIMATEDATA[6][74] = 2026.071$ است.

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

بنابراین، باید سوگیری‌ها را از $BIASDATA[6]$ که در همان موقعیت‌های ۷۳ و ۷۴ اندیس شده‌اند، بازیابی کنیم که عبارتند از $BIASDATA[6][73] = 134.1804$ و $BIASDATA[6][74] = 131.071$.
تخمین صحیح در بازه زیر است:

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

و برای محاسبه تقریب تصحیح شده، می‌توانیم آن مقادیر را درون‌یابی کنیم، به عنوان مثال با استفاده از جستجوی k -نزدیک‌ترین همسایه یا فقط با یک درون‌یابی خطی $y(x) = ax + b$ ، که در آن $y(2003, 1804) = 1869, 0$ و $y(2026, 071) = 1895, 0$.
بنابراین، با استفاده از محاسبات ساده، خط درون‌یابی است

$$y = 1.135837x - 406.28725$$

و مقدار درون‌یابی شده برای تخمین تعداد منحصربه‌فرد ما است

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

با توجه به آزمایش‌های نویسندگان $++HyperLogLog$ ، تخمین n_{lin} که بر اساس الگوریتم شمارش خطی ساخته شده است، حتی در مقایسه با مقدار تصحیح شده با سوگیری n برای تعداد منحصربه‌فردهای کوچک هنوز بهتر است. بنابراین، اگر حداقل یک شمارنده خالی وجود داشته باشد، الگوریتم علاوه بر این تخمین خطی را محاسبه می‌کند و از فهرستی از آستانه‌های تجربی، که در جدول ۷ می‌توان یافت، برای انتخاب ارزیابی ترجیحی استفاده می‌کند. در چنین حالتی، مقدار تصحیح شده با سوگیری n تنها زمانی استفاده می‌شود که تخمین خطی n_{lin} بالاتر از آستانه χ_m برای m فعلی قرار گیرد.

مثال ۱۱- تصحیح سوگیری با آستانه

مثال ۱۰ را در نظر می‌گیریم، که در آن برای $m = 2^{10}$ مقدار تصحیح شده با سوگیری $n \approx 1886.218$ را محاسبه کردیم. برای تعیین اینکه آیا باید این مقدار را به تخمین شمارش خطی ترجیح دهیم یا خیر، باید تعداد شمارنده‌های خالی Z را در داده‌ساختار ابرلگ‌لگ++ پیدا کنیم. از آنجایی که آن مقدار را در مثال خود نداریم، فرض می‌کنیم مقدار مذکور برابر $Z = 73$ است. بنابراین، تخمین خطی طبق معادله تخمین برابر است با:

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

سپس، n_{lin} را با آستانه $\chi_m = 900$ از جدول ۷ مقایسه می‌کنیم که بسیار کمتر از مقدار محاسبه شده است، بنابراین تخمین تصحیح شده با سوگیری n را به تخمین شمارش خطی n_{lin} ترجیح می‌دهیم.

جدول ۷- آستانه‌های تجربی χ_m برای مقادیر دقت پشتیبانی شده

p	m	χ_m	p	m	χ_m	p	m	χ_m
۴	۲ ^۴	۱۰	۹	۲ ^۹	۴۰۰	۱۴	۲ ^{۱۴}	۱۱۵۰۰
۵	۲ ^۵	۲۰	۱۰	۲ ^{۱۰}	۹۰۰	۱۵	۲ ^{۱۵}	۲۰۰۰۰
۶	۲ ^۶	۴۰	۱۱	۲ ^{۱۱}	۱۸۰۰	۱۶	۲ ^{۱۶}	۵۰۰۰۰
۷	۲ ^۷	۸۰	۱۲	۲ ^{۱۲}	۳۱۰۰	۱۷	۲ ^{۱۷}	۱۲۰۰۰۰
۸	۲ ^۸	۲۲۰	۱۳	۲ ^{۱۳}	۶۵۰۰	۱۸	۲ ^{۱۸}	۳۵۰۰۰۰

الگوریتم کامل ابرلگ++ در زیر نشان داده شده است.

الگوریتم ۱۱- تخمین تعداد منحصربه‌فرد با ابرلگ++

ورودی: مجموعه داده D

ورودی: آرایه‌ای از m شمارنده لگ با تابع درهم‌ساز h

خروجی: تخمین تعداد منحصربه‌فرد

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

برای $x \in D$ انجام بده

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

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

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

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

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

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

$$n \leftarrow \hat{n}$$

اگر $\hat{n} \leq 5m$ آنگاه

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

$$Z \leftarrow \sum_{j=0}^{m-1} (COUNTER[j])^2$$

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

ویژگی‌ها

دقت ابرلگ++ در طیف وسیعی از محاسبات تعداد منحصربه‌فردها بهتر از ابرلگ است و برای بقیه به همان اندازه خوب است. برای تعداد منحصربه‌فردهای بین ۱۲۰۰۰ و ۶۱۰۰۰، تصحیح سوگیری امکان خطای کمتری را فراهم می‌کند و از افزایش ناگهانی خطا هنگام جابجایی بین الگوریتم‌های فرعی جلوگیری می‌کند.

با این حال، به دلیل عدم نیاز ابرلگ++ به ذخیره مقادیر درهم‌ساز، فقط یک به اضافه حداکثر اندازه تعداد صفرهای ابتدایی، حافظه موردنیاز در مقایسه با ابرلگ رشد قابل توجهی نمی‌کند و با توجه به معادله تخمین حافظه به 6×2^p بیت نیاز دارد. الگوریتم

ابرلگ++ تخمین تعداد منحصر به فرد حدود 7.9×10^9 مقدار با میزان خطای معمولی ۱.۶۲۵ درصد با استفاده از ۲.۵۶ کیلوبایت حافظه استفاده می شود.

همانطور که قبلا اشاره شد، الگوریتم از رویکرد میانگین گیری تصادفی استفاده می کند و مجموعه داده را به $m = 2^p$ زیرمجموعه $\{S_j\}_{j=1}^{m-1}$ تقسیم می کند که هر کدام دارای شماره های مرتبط $\{COUNTER[j]\}_{j=1}^{m-1}$ هستند، هر شماره اطلاعات مربوط به $\frac{n}{m}$ مقدار را مدیریت می کند. هویله و همکاران متوجه شدند که برای $m \ll n$ بیشتر شماره ها هرگز استفاده نمی شوند و نیازی به ذخیره سازی ندارند، بنابراین ذخیره سازی می تواند از نمایشی تنک بهره گیرد. اگر تعداد منحصر به فرد n بسیار کوچکتر از m باشد، ابرلگ++ به حافظه قابل توجهی کمتری نسبت به روش های قبلی خود نیاز دارد.

الگوریتم ابرلگ++ در نسخه تنک صرفا جفت های $(j, COUNTER[j])$ را ذخیره می کند و آنها را با الحاق الگوهای بیتی خود به عنوان یک عدد صحیح منفرد نشان می دهد. همه این جفت ها در فهرست مرتب شده منفرد از اعداد صحیح ذخیره می شوند. چون همیشه در پی محاسبه حداکثر رتبه هستیم، نیازی به ذخیره جفت های مختلف با اندیس یکسان نیست. در عوض فقط جفتی با حداکثر اندیس باید ذخیره شود. در عمل، برای ارائه نتیجه بهتر، می توان فهرست مرتب نشده دیگر را برای درج های سریع حفظ کرد که باید به طور دوره ای مرتب شده و در فهرست اصلی ادغام شوند. اگر چنین فهرستی به حافظه بیشتری نسبت به نمایش متراکم شماره ها نیاز داشته باشد، می توان به راحتی به شکل متراکم تبدیل شد. علاوه بر این، برای نمایش تنک فشرده تر، الگوریتم ابرلگ++ تکنیک های فشرده سازی مختلفی را با استفاده از رمز گذاری با طول متغیر و رمز گذاری تفاوت برای اعداد صحیح پیشنهاد می کند، بنابراین فقط اولین جفت و تفاوت ها از مقدار آن ذخیره می شوند. در حال حاضر، الگوریتم ابرلگ++ به طور گسترده در بسیاری از برنامه های پر کاربرد، از جمله Google BigQuery و Elasticsearch استفاده می شود.