

Universal Online Sketch for Tracking Heavy Hitters and Estimating Moments of Data Streams

2020 INFOCOM Toronto (Q1)

3 citations

Shigang Chen: University of Florida (Computer Networks, Big Data, Internet Security)

اندازه گیری ترافیک در مسایل مختلف حوزه مانیتورینگ شبکه مثل تشخیص حملات، تشخیص آنتروپی می تواند موثر باشد. برای تشخیص آنومالی (یکی از راهکارهای تشخیص حملات) بسته هایی که به صورت استریم وارد می شوند را بر اساس هدرشان به جریان های مختلفی تقسیم میکنند و حال اطلاعات آماری آن جریانها را به دست می آورند^۱ و از این اطلاعات آماری برای تشخیص بی نظمی استفاده می شود. با توجه به منابع محدود دستگاه های میانی، الگوریتم های زیرخطی^۲ مختلفی به نام اسکچ های عمومی مثل: count min و count sketch ارایه شدند. اما مشکلشان این بود که هر کدام برای یکی از مسایل اندازه گیری طراحی شده بودند و فراهم کردن اسکچ های مختلف برای مسیله های دیگر فضای زیادی می طلبید. لذا Universal sketch هایی بایستی طراحی می شد که نه تنها جریان های بزرگ heavy hitter را تشخیص می دادند بلکه اطلاعات آماری کلی نیز در مورد توزیع ترافیک یا همان moment ها نیز ارایه می دادند. یکی از روش های موفق در این زمینه univmon می باشد که دراصل چندین جدول هشینگ دارد که هر بسته معمولاً چندبار ذخیره^۳ می شد و لذا حجم زیادی را در بر می گرفتند و سرعت مان را نیز کند می کردند. لذا قصد داریم روشی به نام LUV ارایه دهیم که از اسکچی بر پایه نمونه برداری تدریجی^۴ (به جای نمونه برداری سلسه مراتبی^۵) می باشد که از ActiveCM+ در عوض CS استفاده می کند که سربار عملیات هر بسته را می کاهد و دارای نرخ گذر بالا (مانند اسکچ های تکی) می باشد. در خلاصه باید گفت، اسکچ های قبلی به منظور اندازه گیری جریان می بودند و اسکچ های جهانی نیز مشکلات سربار دسترسی حافظه دارند.

مانیتورینگ شبکه با مفاهیم زیر مرتبط می باشد:

۱- Flow moment: اطلاعاتی که از تمامی جریان های استریم در یک لحظه می توان به دست آورد. در واقع وضعیت کل ترافیک را نشان می دهد. می توان در بازه های زمانی lg را حساب کرد و به اصطلاح یک سری زمانی از اون معیار به دست آورد و سپس با مقایسه آن با یک مقدار آستانه ای، رویدادهای آنومالی را تشخیص دهیم.

$$L_g = \sum_{1 \leq f \leq F} g(n_f)$$

این که تابع g چی باشد، شناسه های مختلفی از ترافیک را نشان می دهد که می تواند در عملیات مانیتورینگ مفید واقع شود. (به مقاله مراجعه شود). برای مثال برای شناسایی حملات منع خدمت می توان هاست هایی که به سمت آنها جریانات بیش از حدی از طرف منابع مختلف می باشد را شناسایی کنیم. برای اینکار می توان از یک اسکچ برای ذخیره سازی هاست های موردنظر و $g(x)=x^0$ قرار داد.

۲- Heavy hitter: جریانی است که سائز آن تاثیر زیادی بر روی یک flow moment می گذارد یا به عبارت ساده حجم زیادی از ترافیک را تشکیل می دهند:

$$H_g = \{f \mid g(n_f) \geq \alpha L_g\}.$$

¹ Flow-level statistics

² Sub-linear

³ record

⁴ Progressive sampling

⁵ Hierarchical sampling

ضبط هردوی heavyhitter ها و moment ها به شناسایی آنومالی‌ها و مانیتور شبکه کمک خواهد کرد. به عنوان مثال برای تخمین آنتروپی و کنترل ترافیک ارسالی به یک IP مشخص، می‌توان جریانات را بر اساس آدرس مبدا تقسیم کرد و آنتروپی جریان را برابر مقدار زیر تعریف کرد. حالا هر تغییری شگرفی در این مقدار به معنای تحت حمله قرار گرفتن یک آدرس می‌باشد.

$$E = - \sum_{1 \leq f \leq F} \frac{n_f}{n} \log \frac{n_f}{n}$$

معیارهایی که برای ما مهم است:

- Heavy hitter estimation
- Moment Estimation
- Memory Overhead: اسکچ‌ها چون روی سوییچ‌ها و روترها می‌باشند، لذا فضای حافظه کمی دارند و باید بین دقت اسکچ و میزان فضای موردنیاز، سبک-سنگین^۶ کرد. باید اون کران پایین دقت محاسبات را بتوان تضمین کرد.
- Packet-Processing cost: محاسبات هشینگ و دسترسی به حافظه که برای هر بسته صورت می‌گیرد. هشینگ نیز به دلیل نمونه برداری بسته یا برای دستیابی به مکان تصادفی در حافظه می‌باشد. از آنجا که برای هر بسته نمونه برداری تنها یک بار انجام می‌شود، بیشتر هزینه به خاطر محاسبات هشینگ متعدد برای نوشتن/خواندن حافظه می‌باشد.

راهکارهای آرایه شده قبلی:

- اسکچ‌های اندازه‌گیری سایز هر جریان: Count Min از یک آرایه دوبعدی شمارنده‌ها شامل d سطر که معمولاً برابر ۴ می‌باشد. هر بسته را d بار به یکی از این شمارنده‌ها هش می‌کند و مقدارش را افزایش می‌دهد. Conservative Update همه هش‌ها را می‌خواند و کوچکترین آنها را آپدیت می‌کند. اگر سایز یک فلو را کویری بزینم، این دو اسکچ کوچکترین مقدار را از بین d تا هش انتخاب می‌کنند و تخمین می‌زنند. Count Sketch به جای افزایش شمارنده‌های هش شده، مقدار +1/-1 هش شده به شمارنده‌های هش شده اضافه می‌کند. که تابع هش +1/-1 به صورت شبه تصادفی، یک جریان را به +1 یا -1 نگاشت می‌کند. و وقتی کویری بزینم مقدار میانگین یا میانه اون d تا شمارنده را برمی‌گرداند. VAC و.. نیز از جمله راهکارهای دیگر می‌باشند.
- برای تمامی اسکچ‌های گفته شده، دقت اندازه‌گیری به دو عامل بستگی دارد:
 1. اگر تعداد بسته‌های جریان‌ها افزایش یابد، شمارنده‌های اشتراکی جریان‌ها (مثلاً به دلیل تصادم هش) نرخ خطای بالا خواهند داشت.
 2. حافظه اختصاص داده‌شده به اسکچ هر چه بیشتر باشد، شمارنده‌ها افزایش می‌یابد و نرخ خطا به دلیل اشتراک کمتر، کاهش خواهد یافت
- اسکچ‌های جهانی: توانایی اجرای چندین عملیات اندازه‌گیری مختلف را دارند. دیگر ماژول‌های جدا جدا برای اندازه‌گیری سایز جریان، شناسایی heavy hitter ها و moment ها لازم نیست. OpenSketch از یک سری ماژول‌هایی برای استفاده مشترک بین چندین task، استفاده می‌کند. اما Univmon از راهکار دیگری استفاده می‌کند. از چندین CS (رسماً ۱۵ تا) استفاده می‌کند. با احتمال $1/2^i$ جریان‌ها را نمونه برداری و پکت‌هایشان را در CS مربوطه ضبط می‌کند.

حال به معرفی نوآوری‌های استفاده شده در LUV می‌پردازیم:

نمونه‌برداری تدریجی: به جای نمونه برداری سلسله مراتبی که در univmon استفاده می‌شود و هر بسته را در تعداد متغیری اسکچ (میانگین ۴ تا) ضبط می‌کند، هر بسته را دقیقاً در یک اسکچ ذخیره می‌کند لذا شمارنده‌ها نویز کمتری خواهند داشت و دقت بالاتر خواهد رفت. روش ما و univmon هر دو از چندین اسکچ برای اندازه‌گیری استفاده می‌کنند.

⁶ tradeoff

ActiveCM+ Sketch: حافظه کمتر، هش‌های کمتر و بارحافظه کمتر نسبت به univmon استفاده می‌کند. LUS برای هر بسته حداکثر ۵ بار هش و ۲ بار دسترسی به حافظه نیاز خواهد داشت. دلیل دقت بهتر در حافظه یکسان نسبت به CS استفاده از شمارنده‌های بیشتر می‌باشد، چون که خروجی شمارنده‌ها نیز تعداد بیت کمتری خواهند داشت لذا می‌توان تعداد شمارنده‌ها را افزایش داد. داده ساختار فشرده ActiveCM Sketch، داده‌ها را به صورت احتمالی ضبط می‌کند و نسخه بهبود یافته آن ActiveCM+، امکان شناسایی heavy-hitter ها را نیز خواهد داشت. ActiveCM از یک داده ساختار فشرده و دو روش زیر استفاده می‌کند:

- Counter sharing
- Counter compression

ActiveCM از عملیات زیر پشتیبانی می‌کند:

- ضبط بسته رسیده شده: زمانی که بسته جدید رسید، تمامی شمارنده‌های مجازی متناظر با جریان را یک واحد افزایش می‌دهد (البته با احتمال $\frac{1}{2^{PCF[i]-a}}$) که البته می‌توان آن را به صورتی تغییر داد که تعداد بایت‌های جریان را نیز بشمارد. همانطور که می‌دانیم d تا هش به ازای هر بسته بایستی اجرا کرد. اما به جای اینکار، یک هش ۶۴ بیتی یکجا از جریان حساب می‌کنیم و اونو به d قسمت تقسیم می‌کنیم.
- کویری و برگرداندن سائز یک جریان: کوچکترین مقدار از بین d شمارنده جریان را برمیگرداند.

ActiveCM+ : سائز هر جریان را تخمین می‌زند مثل CM و CS اما با حافظه‌ای کمتر. بعد از ارتقاش میدیم با استفاده از یک مین هیپ k ، تا جریان و شمارنده متناظرشان را شناسایی کند ← heavy hitter. همانند یک فیلتر اولیه از هیپ کمینه استفاده می‌کند.

کران‌دار بودن: در این بخش ثابت می‌کند که مقدار کویری بازگردانده شده، کران دار می‌باشد:

$$Pr\{\hat{n}_f \geq n_f + (n - n_f) \frac{2d}{m}\} \leq \left(\frac{1}{2}\right)^d$$

اجرا و ارزیابی: در یک حجم محدودی از حافظه، ActiveCM+ را با اسکچ‌های دیگر شناسایی heavy hitter ها از نظر کارایی مقایسه می‌کنیم. سپس LUV را با univmon برای تخمین در لحظه moment ها مقایسه می‌کند. نرخ گذر بسته‌ها که متاثر از سرعت عملیات هش (درج بسته‌ها) و دسترسی به حافظه (درج و کویری) هست، را نیز بررسی می‌کنیم. داده‌های آزمایش نیز از دیتاست cadia می‌باشند.

- شناسایی Heavy hitter ها با یک حافظه محدود :

- نرخ گذر: در مقایسه با اسکچ‌های دیگر نرخ گذر بالاتری خواهد داشت.
- دقت تخمین

▪ Average Relative Error(ARE)

- میزان تاثیر اندازه جدول مجازی یا همان d را بر روی دقت ActiveCM+ نیز بررسی می‌کنیم. و مقدار بهینه آن را به دست می‌آوریم

- مقایسه با Universal Sketch های دیگر: Univmon از ۱۴ لایه سلسله مراتبی استفاده می‌کند. درحالی که LUS از ۱۰ تا subsketch که ActiveCM+ شامل ۴ ستون می‌باشند، استفاده می‌کند. در موارد زیر برتر می‌باشد:

- نرخ گذر بسته‌ها
- دقت تخمین سائز جریان. چون بسته‌ها را تنها در یک subsketch نمونه‌برداری شده، وارد می‌کند، نویز و خطای کمتری خواهد داشت.

- همچنین تاثیر تعداد subsketch را نیز بررسی کردیم. اگر در univmon لایه‌ها را افزایش دهیم، دقت افزایش می‌یابد اما در LUS با افزایش تعداد subsketch ها چون آخرین subsketch یک جدول هش می‌باشد، که تمامی جریان‌های نمونه برداری شده را در خود جای داده است، این امر تاثیری نخواهد داشت.