

Sketch-based Change Detection: Methods, Evaluation, and Applications

ACM 2003

B Krishnamurthy (AT&T), Yin Zhang (University of Texas Austin), ...

742 citations

Related Keywords: change detection, DataStream Computation, Sketch

تشخیص بی‌نظمی‌ها^۱ (تغییرات) در شبکه‌های پویای امروزی امری مهم می‌باشد. برخی از این بی‌نظمی‌ها قابل پیش‌بینی، عادی و تحمل‌پذیر می‌باشند، اما برخی دیگر نشانی از یک حمله می‌توانند باشند. دو روش شناسایی بی‌نظمی وجود دارد: مبتنی بر امضا^۲ یا روش‌های آماری^۳ که امکان شناسایی حملات جدید را نیز دارا می‌باشد و یک مدلی از رفتار نرمال شبکه به دست می‌آورد. الگوریتم‌های تشخیص بی‌نظمی، معمولاً ترافیک را به صورت مجموعه‌ای از جریان‌ها می‌بینند و سعی می‌کنند تغییرات را تشخیص دهند^۴. اما با افزایش حجم داده‌ها، نگهداری اطلاعات هر جریان^۵ ناممکن است. اما ما در این مقاله از داده ساختاری جدید بر مبنای اسکچ‌ها^۶ که اطلاعات را به صورت فشرده نگهداری می‌کنند، به نام k-array استفاده می‌کنیم. از میزان ثابت و ناچیزی از حافظه استفاده می‌کند و هزینه به روزرسانی و بازسازی هر رکورد نیز ثابت است. سپس یکسری مدل‌های پیش‌بینی سری زمانی^۷ (ARIMA, HoltWinters, ...) طبق این اطلاعات فشرده پیاده‌سازی می‌کنیم و جریان‌هایی که بیشترین خطا از مقدار پیش‌بینی شده را داشتند، شناسایی می‌کنیم. همچنین یک سری مکاشفه برای تنظیم پارامترهای خودکار مدل نیز ارائه می‌دهیم. سپس با استفاده از داده‌های جمع‌آوری شده از یک ISP، نتیجه می‌گیریم که روش ارائه شده دقت بالایی دارد و می‌تواند در دستگاه‌های با سخت‌افزار پایین پیاده‌سازی شود.

روش‌های قدیمی تشخیص تغییر بر پایه پیش‌بینی سری زمانی و آنالیز داده‌های پرت در نرخ ترافیک‌های حجیم امروزی با تعداد زیادی از سری‌های زمانی جواب نمی‌دهند چون برای هر جریان می‌خواهند اطلاعاتی نگهداری کنند. روشی که می‌خواهیم ارائه دهیم در زمینه data stream computation (اون دو ویژگی معروف یعنی بررسی هر تاپل تنها یک بار و یا به مقدار کم و استفاده محدود از حافظه را دارا می‌باشد) می‌باشد. اسکچ‌ها یکی از تکنیک‌های موجود در این زمینه می‌باشند که با استفاده بهینه از حافظه، از ویژگی‌های تضمین بازسازی با نرخ احتمال بالای قابل اثبات و خطی بودن برخوردار می‌باشند.

مدل‌های دیتاستریم: مدل‌های مختلفی برای توصیف دیتاستریم‌ها ارائه شده است: Time Series Model, Cache Register Model, Turnstile Model. که در اینجا از Turnstile Model استفاده می‌کنیم:

¹ Anomaly

² Signature-Based

³ Statistical Approaches

⁴ Change Detection

⁵ Per-flow analysis

⁶ Sketch

⁷ Time-series forecasting

مقادیر مختلف کلید $[u] = \text{key space}$ I یک ست استریم شامل توالت‌های زیر است.

$$I = \alpha_1, \alpha_2, \dots$$

$$\alpha_i = \{(a_i, v_i) \mid a_i \in \{1, 2, \dots, u-1\}, v_i \in \mathbb{R}\}$$

هر کلید a_i مقدار $A[a_i]$ دارد که هرگاه یک بایت جدید (a_i, v_i) می‌آید، $A[a_i] \leftarrow v_i$ **که سیگنال**

هدف این است که سیگنال‌های با تغییرات شاخص را شناسایی کنیم. مدل ارایه شده بسیار کلی است و می‌تواند در زمینه‌های مختلف به کار رود: مثلاً کلید می‌تواند همان مشخصات flow ID باشد و آپدیت هم اندازه ی جریان (به بایت) باشد. در این مقاله ما از تنها آدرس مقصد و جمع بایت‌ها در آزمایش‌هایمان استفاده کرده‌ایم.

معماری:

روش ما شامل سه کامپوننت زیر است:

- Sketch module
- Forecasting Module
- Change Detection Module

برای یک *substream* از $\{a_1, a_2, a_3, \dots\}$ **ست** $(a_1, v_1), (a_2, v_2), (a_3, v_3), \dots$ **Sketch module**

Second moment

$$\forall a \in [u]: \gamma_a = \sum_{i \in A_a} v_i^2, A_a = \{i \mid a_i = a\}$$

$$F_2 = \sum_a \gamma_a \rightarrow \text{L2 norm} = \sqrt{F_2}$$

اسکچ‌ها از یک داده ساختار فشرده برای ذخیره \mathcal{V}_a در هر دوره زمانی می‌کند. اسکچ‌ها داده ساختارهای احتمالی بر مبنای random projection می‌باشند. یک اسکچ خاص به نام k-array ارزیابی می‌دهیم که بر پایه count sketch می‌باشد. اما عملیات روی آن ساده تر و کاراتر می‌باشد. یک اسکچ S از $H \times K$ جدول رجیستری تشکیل شده است:

$$T_S[i][z] : i \in [H], z \in [K]$$

سایز جدول K و تعداد توابع هشی H

هر سطر یک تابع هش $h_i: [u] \rightarrow [k]$ می باشد. در واقع این داده ساختار، آرایه ای از جداول هش می باشد.
توابع هش بایستی 4-universal باشند تا دقت بازسازی تضمین شود. عملیاتی که داریم:

- Update : به روزرسانی اسکچ
- Estimate : بازسازی v_a برای کلید a
- estimateF2
- combine : ترکیب خطی چندین اسکچ

- Update(S, a, v): $\forall i \in [H], T_S[i][h_i(a)]_+ = v$

- Estimate(S, a): *rest*

$$v_a^{\text{estimated}} = \text{median}_{i \in [H]} v_a^{h_i}$$

$$v_a^{h_i} = \frac{T[i][h_i(a)] - \frac{\text{sum}(S)}{K}}{1 - \frac{1}{K}}$$

$$\text{sum}(S) = \sum_{j \in [K]} T_S[i][j]$$

- Estimate $F_2(S)$:

$$F_2^{h_i} = \frac{K}{K-1} \sum_{j \in [K]} (T_S[i][j])^2 - \frac{1}{K-1} (\text{sum}(S))^2$$

$$F_2^{\text{est}} = \text{median}_{i \in [H]} F_2^{h_i}$$

- Combine($c_1, S_1, \dots, c_L, S_L$):

$$S = \sum_{k=1}^L c_k \cdot S_k, \quad T_S[i][j] = \sum_{k=1}^L c_k \cdot T_{S_k}[i][j]$$

Forecasting module: از اسکچه‌های مشاهده شده در بازه های قبلی $S_0(t' < t)$ برای ساخت اسکچ پیشبینی $S_f(t)$ و سپس محاسبه خطای $Se(t)$ استفاده می‌کند. مدل برای پیش بینی سری زمانی تک متغیره و تشخیص تغییر استفاده می‌کنیم. چهار تایی اول مدل‌های نرم کننده⁸ ساده و بقیه از خانواده ARIMA می‌باشند.

Moving Average:

$$S_f(t) = \frac{\sum_{i=1}^w S_f(t-i)}{w} \quad (w > 1)$$

وزن‌ها به هم نمونه‌های قبلی و وزن یکسانی دهنده

پارامترها w : تعداد نمونه‌های قبلی

مدل‌های Arima: یک دسته از تکنیک‌های پیشبینی سری‌های زمانی می‌باشند. که می‌توانند رفتارهای متفاوتی را برای پیش‌بینی سری‌های زمانی تک متغیره و تشخیص تغییر مدل کنند.

یک مدل $ARIMA(p, d, q)$

معنی

p : # autoregressive parameters

q : # moving average parameters

d : # differencing passes

پارامترهای بالایی پس از پذیرش تفاوت محاسبه می‌شوند.

$$z_t - \sum_{i=1}^q MA_i \cdot z_{t-i} = C + e_t - \sum_{j=1}^p AR_j \cdot e_{t-j}$$

z_t : difference کردن سری زمانی اصلی

e_t : اوردوینگ بین در زمان t

MA_i, AR_j ($1 \leq i \leq q, 1 \leq j \leq p$) ضرایب ثابت اند.

ARIMA0: ($p \leq 2, d=0, q \leq 2$)

ARIMA1: ($p \leq 2, d=1, q \leq 2$)

درجته اسکچها دونه ARIMA در تری گریم:

Change detection module: بعد از محاسبه $S_e(t) = S_0(t) - S_f(t)$ این ماژول بر اساس مقدار T_A تصمیم می‌گیرد.

$$T_A = T \cdot \left[\text{ESTIMATEF2}(S_e(t)) \right]^{\frac{1}{2}}$$

برای هر کلید حال کافی است $\text{Estimate}(S_e(t), a)$ را با T_A مقایسه کنیم.

این که به چه نحوی مقادیر کلید را ذخیره کنیم تا بتوانیم آنها را به این ماژول بدهیم، یک مساله است. روشی که استفاده می‌کنیم (به جای نگهداری اطلاعات هر جریان) $S_e(t)$ را در ابتدا محاسبه می‌کنیم و سپس تغییرات را در مرحله بعد حساب می‌کنیم. چون خود استریم ورودی تمام کلیدها را دارا می‌باشد نیازی به ذخیره همه آنها برای هر جریان نیست. اما نیاز به دسترسی مجدد به استریم دارد و لذا در موارد آفلاین کاربرد دارد.

تنظیم پارامترها:

پارامترهای H و K اسکیج می‌تواند در کارایی و مطمئن بودن ارورهای پیش‌بینی تاثیرگذار باشد. یکی از توابع مکاشفه ای که می‌توانیم استفاده کنیم (برای پارامترهای مدل)، کمینه کردن $\sum_t F_2^{est}(S_e(t))$ (انرژی اضافی) می‌باشد. الگوریتمی که برای این کار استفاده می‌کند نیز جستجوی multi-pass grid می‌باشد (مثلاً برای مدل EMWA).

برای پارامترهای مدل: به مقاله مراجعه کن.

راه اندازی محیط آزمایشگاهی: دیتاست شامل داده‌های نت فلو که از روترهای موجود در یک ISP گرفته شده است.

نتایج آزمایش: پس از تعیین این که چه پارامترهایی برای هر مدل مناسب است (بر اساس grid search)، باید درستی و دقت اسکیج‌ها را در مقایسه با per-flow ها نیز بررسی کنیم. برای $K=32k$ نتیجه می‌شود که مثبت کاذب⁹ و منفی کاذب¹⁰ بسیار کاهش می‌یابد.

کارهای آینده:

⁹ False positive

¹⁰ False negative

- شناسایی تغییرات به خط^{۱۱}: برای شناسایی تغییرات به صورت بلادرنگ^{۱۲} یکی از الزامات اینست که پارامترهای مدل پیش‌بینی به کمک داده‌های قبلی به صورت دوره‌ای محاسبه شوند تا با تغییرات موجود در رفتار ترافیک سازگار باشند.
- دوری از اثرات مرزی به دلیل سائز *interval* های ثابت: راهکارهای ممکن شامل: اجرای چندین مدل به صورت همزمان با *interval* های مختلف و نقاط شروع متفاوت، یا تصادفی کردن سائز *interval* می‌باشد.
- کاهش نرخ مثبت کاذب: یکی از مشکلاتی که برای ترافیک‌هایی که عادی هستند رخ می‌دهد (بی‌نظمی‌های عادی)، با تنظیم کردن پارامترها با شرایط محیطی مساله می‌توان این را برطرف کرد. یک تکنیک اینست که تنها تغییرات اساسی را گزارش کند.
- کامپایل با نمونه برداری: به همراه اسکچ‌ها می‌تواند مقیاس پذیری را بالا ببرد.
- راهنما برای انتخاب پارامترها: با استفاده از تکنیکی که در مقاله ۵ انتخاب شده می‌توان محدوده شان را تعیین کرد.

¹¹ Online change detection

¹² Real-time