KMVスケッチによる時系列データ要約 と突発検知応用に向けた予備検討

○西川 侑志, Thanapol Phungua-eng, 山本泰生 静岡大学 情報学部

木曽シュミットシンポジウム2024

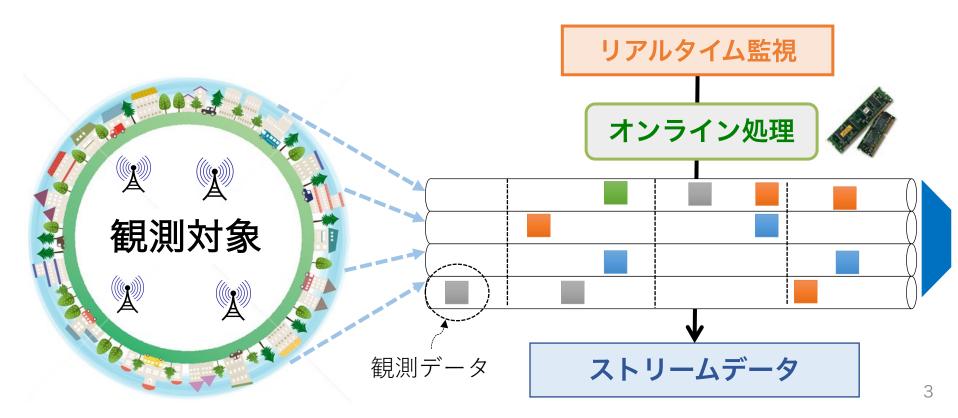
(木曽福島保険センター) 2024年5月15日

発表の概要

- 研究背景
 - ▶ ストリームデータ処理
 - ▶ データ要約
- KMVを用いた頻度サマリの構築
- 頻度サマリを用いた時系列データ要約
 - ▶ 問題設定
 - > アプローチ
- ケーススタディ
 - ▶ 測光時系列データの頻度分布
 - ▶ エビデンスに基づく異常検知

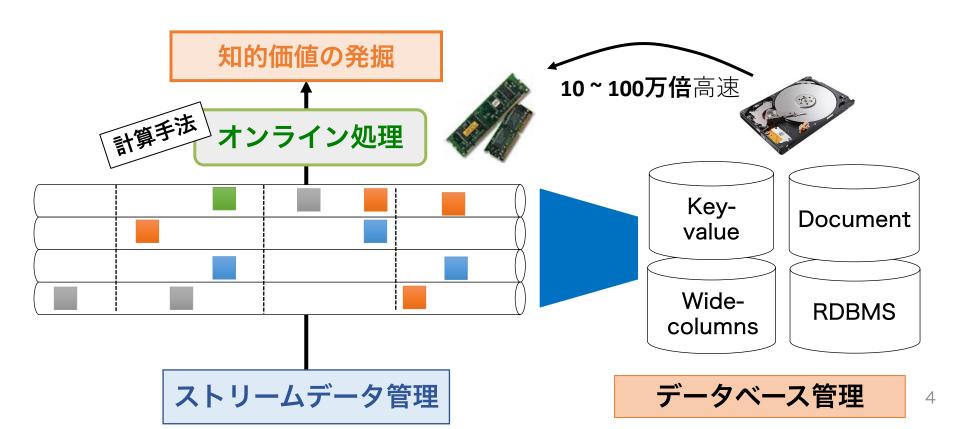
ストリームデータの研究

- ・ストリームデータとは?
 - ▶ 高速に流れ続ける無限長のデータ列
 - ▶ センサーノードから常時到着する観測データ
 - ▶ 観測対象のリアルタイム分析 (傾向の変化や異常の検出)



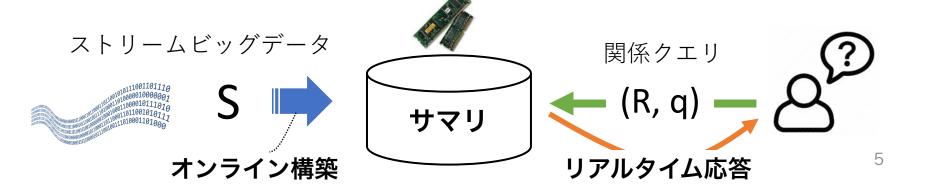
ストリームデータ分析

- ・ビッグデータ分野の重要課題 [IDC, 2015; 2018]
 - ➤ ハードディスク (+SSD) スキャンは原理的に困難
 - ➤ 到着データを "On-the-fly" でインメモリ処理する
 - ▶ 省メモリなインメモリ管理技術の開発が必要不可欠!



データ要約とは?

- ・インメモリ処理を可能にするデータ管理技術
- サマリ: 特殊な関係クエリ (質問) に応答するデータ構造
 - ightharpoonup 管理対象のデータ列: $S = \langle e_1, e_2, ..., e_n \rangle$, n はデータ総数
 - ▶ 関係: R → 問題に応じて設定
- 2種類の関係クエリ(質問)
 - ▶ メンバーシップクエリ: q と関係 R を満たす e_i が S 中に存在した?
 - ▶ サポートクエリ: q と関係 R を満たす e_i が S 中に何回出現した?



発表の概要

- 研究背景
 - > ストリームデータ処理
 - > データ要約
- KMVを用いた頻度サマリの構築
- 頻度サマリを用いた時系列データ要約
 - ▶ 問題設定
 - > アプローチ
- ケーススタディ
 - ▶ 測光時系列データの頻度分布
 - > エビデンスに基づく異常検知

頻度サマリとは

- バイナリ (属性値がOもしくは1) データ集合の要約表現
- 頻度サマリのタスク
 - ▶ 入力: クエリ (属性の集合)
 - ▶ 出力: クエリの頻度

属性数 m

時刻 t	属性 a ₁	属性 a ₂	属性 a ₃	 属性 a _m
1	1	0	1	 0
2	1	1	0	 1
3	1	1	1	 1
· ·				
n	1	0	0	 0

バイナリデータ (a.k.a.トランザクション)の集合

クエリの各属性値が1となるデータ数

例. $q = \{a_1, a_3\}$ のとき, $a_1 = 1$ and $a_3 = 1$ となるデータ数が q の頻度となる



クエリ q

頻度 f(q)

線形探索の場合

時間計算量 O(n)

(データ数に比例して実行時間が増加する)

頻度サマリのアイデア

- 1. バイナリデータの垂直配置
- 2. 各属性の出現時刻集合をスケッチ化

KMVに基づく頻度サマリの構築 (1/2)

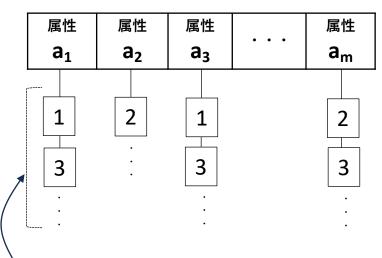
• 頻度サマリのアイデア (その1)

属性毎に値が1となるデータ出現時刻を保持する (バイナリデータ集合の垂直配置)

元のデータ構造

時刻 t	属性 a ₁	属性 a 2	属性 a ₃	 属性 a _m
1	1	0	1	 0
2	0	1	0	 1
3	1	0	1	 1
· ·	•	•	•	 •
n	0	0	0	 0

垂直配置



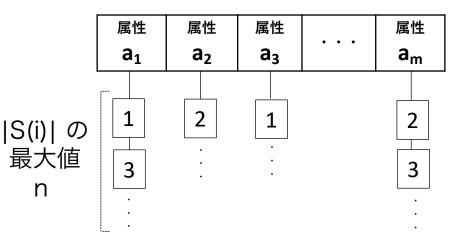
`属性 a₁ が 1となるデータ出現時刻の 集合 {1, 3, …} に相当する

(非出現時刻は記録しない)

KMVに基づく頻度サマリの構築 (2/2)

• 頻度サマリのアイデア (その2)

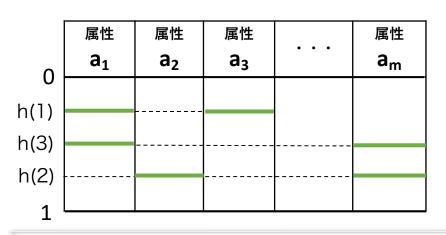
各属性の出現時刻集合を KMV を用いてスケッチ化

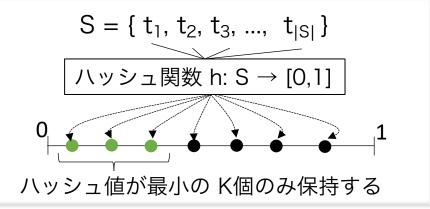


各属性 a_i の出現時刻集合を S_i と書く S_i には多数の時刻が保持される



K-Minimal Values (KMV) を用いて 各 S_i を管理し, 要約





KMVに基づく頻度推定

• S_i : 属性 a_i の出現時刻集合

• *L_i* : *S_i* の KMV (要素数K)

• *U_i* : *L_i* の最大値

• q : クエリ $\{a_{q,1}, a_{q,2}, ..., a_{q,r}\}$

• $\hat{f}(q)$: q の頻度推定値

推定式 [Wang et al., 2021]

$$\hat{f}(q) \Leftrightarrow \frac{K_{\cap}}{\widecheck{U}}$$
 , where

$$K_{\cap} = |L_{q,_1} \cap L_{q,_2} \cap \cdots \cap L_{q,_r}|,$$

$$\check{U} = \max(U_{q,_1}, U_{q,_2}, \cdots, U_{q,_r})$$

 $\hat{f}(q)$ の期待値は f(q) と一致する

ハッシュ関数とKMVの要素数 (K) を増やせば誤差を平均的に小さくできる

元データ

t	a ₁	a ₂	a_3	a ₄
1	1	0	1	1
2	0	1	0	0
3	1	1	1	1
4	1	1	0	0
5	1	0	1	0

頻度サマリ

a_1	a ₂	a ₃	a ₄	
0.4	0.2	0.4 0.3	0.4 0.3	
0.5	0.5	0.7		, '

$$\hat{f}(q) = \frac{1}{0.4} = 2.5$$

$$f(q) =$$

$$h(4) = 0.5$$

 $h(5) = 0.7$

発表の概要

- 研究背景
 - > ストリームデータ処理
 - > データ要約
- KMVを用いた頻度サマリの構築
- 頻度サマリを用いた時系列データ要約
 - ▶ 問題設定
 - > アプローチ
- ケーススタディ
 - ▶ 測光時系列データの頻度分布
 - ▶ エビデンスに基づく異常検知

時系列データ要約とは

- 時系列データの集合を対象とする頻度サマリ
- 頻度サマリのタスク
 - ▶ 入力: クエリ (時系列) + 測定誤差
 - ▶ 出力: クエリの頻度、

時系列データの集合 S = <e₁, e₂, e₃>

クエリの各値 ± 測定誤差の範囲内にある時系列データの個数

時系列 e₂ 観測値 時系列 e₃ 時刻 1 時刻 2 時刻 3 時刻 4

クエリ (時系列) q + 測定誤差 r



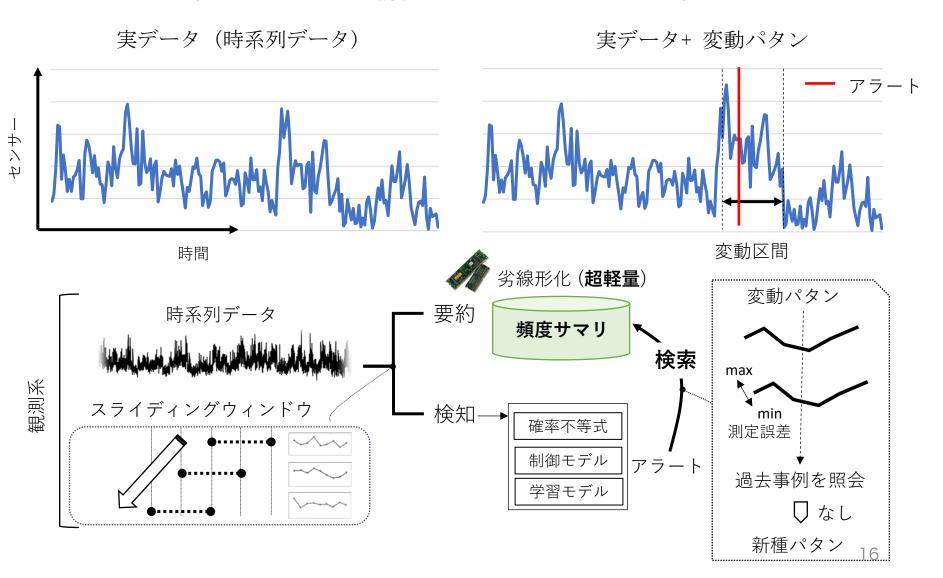
q±rに含まれる時系列は何回出現?



定義 e は q \pm r に含まれる \leftrightarrow $\forall t (q(t) - r \le e(t) \le q(t) + r)$

測光時系列データへの応用可能性

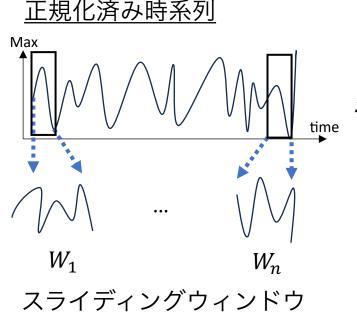
• これまでに見たことのない**新種パタン**をリアルタイム検知する問題



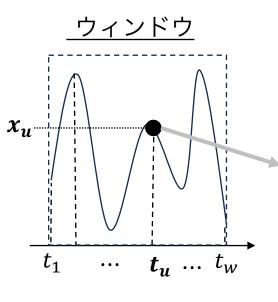
頻度サマリを用いた時系列データ要約

アプローチ

- 1. 時系列データをバイナリデータに変換
 - ➤ 時系列の各値を正規化 (例. Min-Max法, SAX)
 - ▶ 各値の上限と下限をバイナリとして符号化
- 2. バイナリデータの集合から頻度サマリを構築



による時系列切り出し



ウィンドウ幅 = W

x_u	の上限を符号の	Ł
-------	---------	---

0	1	2	3	4	5	6	7	8
0	0	0	0	0	1	1	1	1

 x_{ij} の下限を符号化

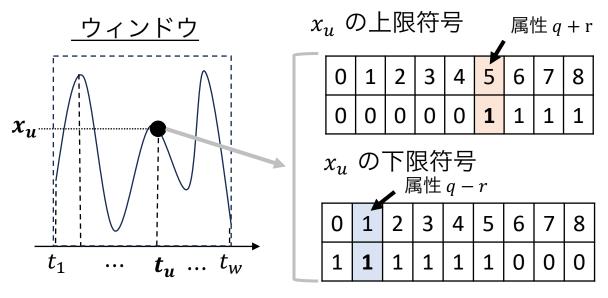
0	1	2	3	4	5	6	7	8
1	1	1	1	1	1	0	0	0

最大値 (Max) = 8 $x_{ij} = 5$ の場合 17

頻度サマリを用いた時系列データ要約

アプローチ

- 1. 時系列データをバイナリデータに変換
 - ➤ 時系列の各値を正規化 (例. Min-Max法, SAX)
 - ▶ 各値の上限と下限をバイナリとして符号化
- 2. バイナリデータの集合から頻度サマリを構築



最大値 (Max) = 8 $x_y = 5$ の場合

上限符号の属性 q+rと 下限符号の属性 q-rの 値が共に1である \leftrightarrow

$$q-r \le x_u \le q+r$$

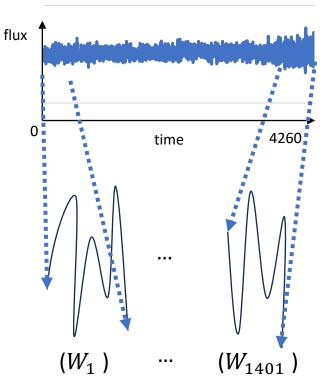
例.
$$q = 3$$
, $r = 2$ のとき $1 \le x_u \le 5$ が成り立つ

発表の概要

- 研究背景
 - > ストリームデータ処理
 - > データ要約
- KMVを用いた頻度サマリの構築
- 頻度サマリを用いた時系列データ要約
 - ▶問題設定
 - > アプローチ
- ケーススタディ
 - ▶ 測光時系列データの頻度分布
 - ➤ エビデンスに基づく異常検知

- 対象ライトカーブ(樫山さんから提供)
 - 108天体*[1]
- ファイルごとに前処理
 - 正規化
 - ➤ Min-Max法:[0~255]に設定
 - ・ウィンドウの切り出し法
 - ▶ウィンドウサイズ:60
 - ▶スライド幅:3





全ウィンドウ数: 1401windows×108天体 = **151,308windows**

例

 $W_1 = \{5032, 5144, 5674, 4867, 4336, \dots, 5232, 4962\}$ $W_2 = \{4867, 4336, 5144, 5313, 4605, \dots, 5213, 5122\}$ $W_{1401} = \{5546, 4365, 5754, 5363, \dots, 5875, 5456\}$

各Windowは シーケンシャル

60timestamps



各値を[0~255]*に正規化

大きな数値を**扱いやすくする**

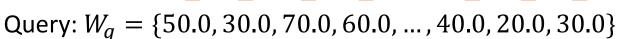
Query(ウィンドウ)

測定誤差: 10%

 $W_u = \{75.5, 55.5, 95.5, 85.5, \dots, 65.5, 45.5, 55.5\}$

(*絶対誤差)







+25.5

 $W_1 = \{24.5, 4.50, 44.5, 34.5, \dots, 14.5, 0.00, 4.50\}$

上限

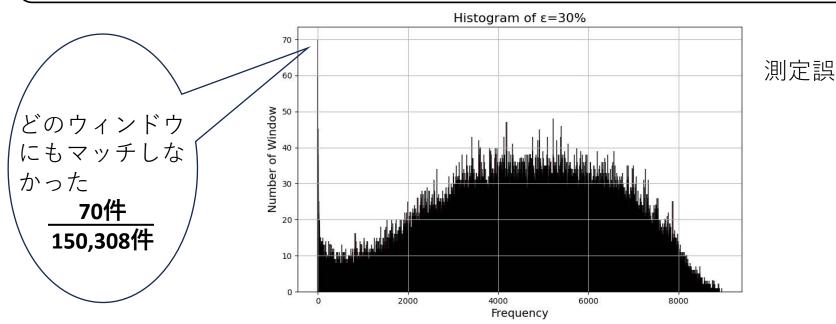
測定誤差: 10% $W_{11} = \{75.5, 55.5, 95.5, 85.5, \dots, 65.5, 45.5, 55.5\}$ (*絶対誤差) +25.5 Query: $W_q = \{50.0, 30.0, 70.0, 60.0, ..., 40.0, 20.0, 30.0\}$ -25.5 $W_l = \{24.5, 4.50, 44.5, 34.5, \dots, 14.5, 0.00, 4.50\}^{\bigoplus l}$ Query W_l (誤差下限) \leq 対象ウィンドウ \leq W_u (誤差上限)

 $W_{1}^{\sim}W_{150,308}$ のうちに測定誤差範囲内のWが何件発生しているのか?

ウィンドウの集合

• 頻度分布: 各ウィンドウの頻度 (真値) 算出

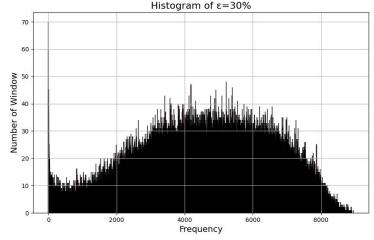
全150,308ウィンドウの中で, 自身以外のウィンドウ(150,307ウィンドウ) に対して, 自身がクエリとして**頻度**を求める



測定誤差: 30%

このような頻度分布のデータセットに対して,調べたいQueryを投げる
→頻度が低ければ,レアな現象(突発現象)

- 予備実験における目的
 - ▶真値による頻度分布をもとに, **突発信号**をクエリと して頻度の算出
- 対象データ
 - ➤ 実フレア (Tomo-e) *[2]
 - ▶逢澤さんから提供
 - ➤ 人工フレア (Keppler) *[3]
 - ➤Thanapolさんから提供



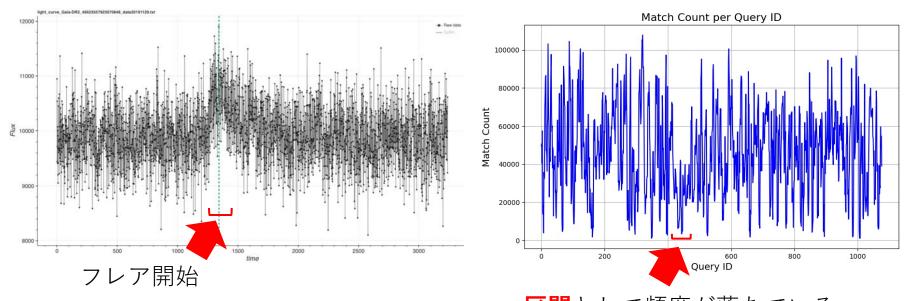
フレアのウィンドウをクエリとして,頻度分布にて比較 →頻度が低ければ,レアな現象(突発現象)

[2] M. Aizawa, K. Kawana, K. Kashiyama, R. Ohsawa, H. Kawahara, F. Naokawa, T. Tajiri, et al. Fast optical flares from M dwarfs detected by a one-second-cadence survey with Tomo-e Gozen. Publications of the Astronomical Society of Japan (PASJ), 74(5):1069–1094, 2022.

・実フレア (Tomo-e)

"light_curve_Gaia-DR2_46623557923070848_date20191129"

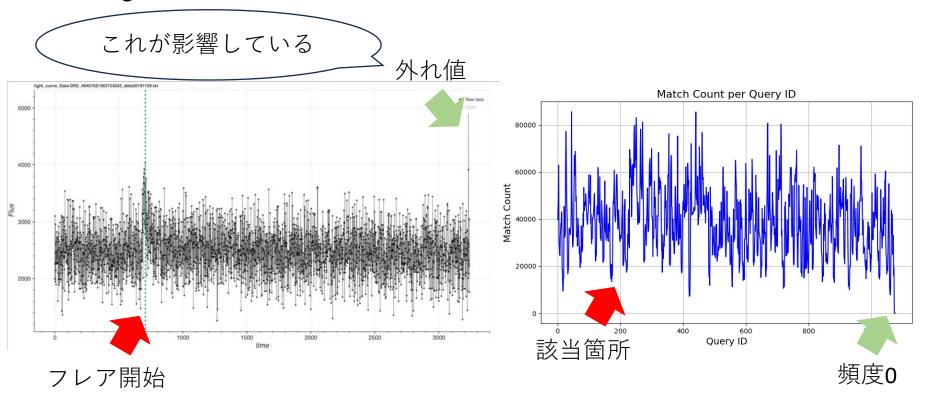
フレア以外で頻度が下がっている原因は,上下が激しいデータだから?



区間として頻度が落ちている

・実フレア (Tomo-e)

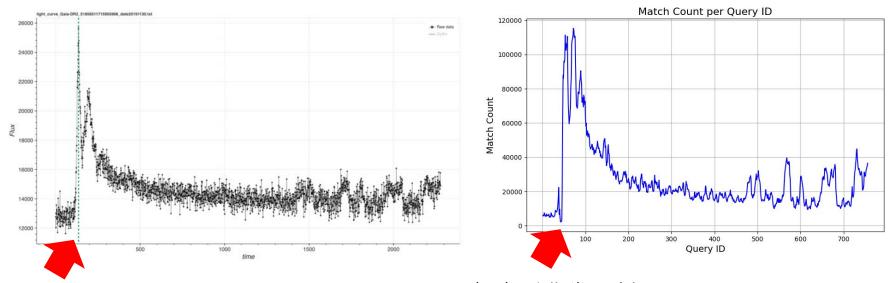
"light_curve_Gaia-DR2_49407521363733632_date20191129"



• 実フレア (Tomo-e)

"light_curve_Gaia-DR2_51856511715955968_date20191130"

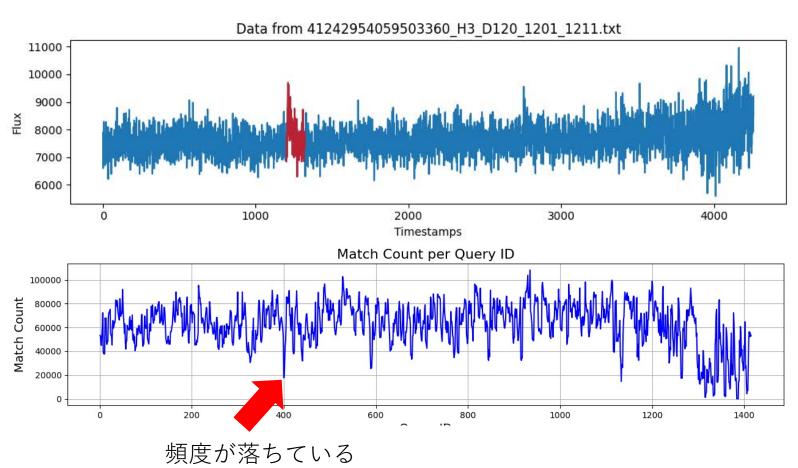
似たようなグラフ結果



フレア開始

頻度が非常に低い

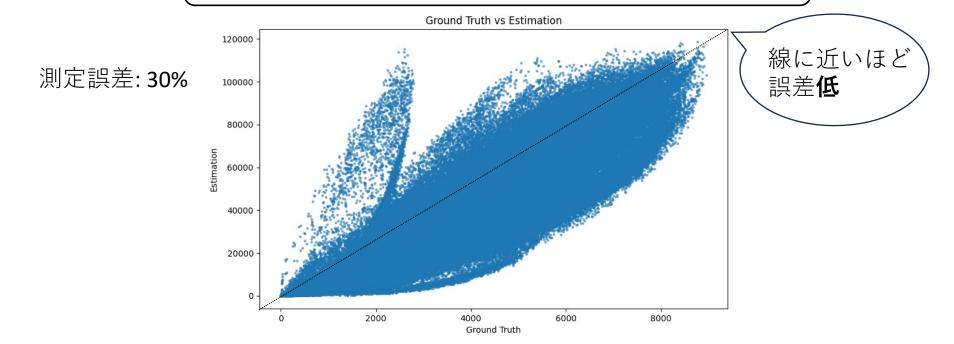
• 人工フレア (Keppler)



・ 真値と推定値の対角グラフ

KMV(データスケッチング)

自身以外の150,107ウィンドウのうち, K=5,000で頻度推定



ハッシュ関数とKMVの要素数 (K) を増やせば, 誤差を平均的に小さくできる

・ 真値と推定値の頻度分布

์KMV(データスケッチング)

自身以外の150,107ウィンドウのうち, K=5,000で頻度推定

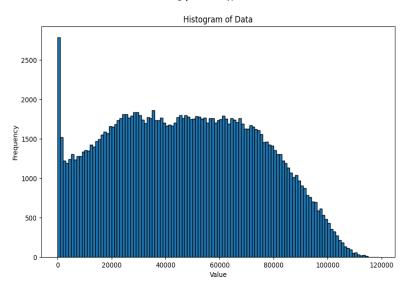


Histogram of ε=30%

Frequency

推定值

測定誤差: 30%



スケールを伸ばしたら,似たような頻度分布 → 少量のデータ(K=5000)で**突発検知の可能性**

まとめと今後の課題

- ・まとめ
 - ▶KMVに基づく頻度サマリの紹介
 - ▶時系列データ要約の提案
 - ▶ライトカーブデータを用いた予備実験
 - 突発クエリの頻度は低くなる
 - 真値と推定値の誤差の確認
- 今後の課題
 - ▶推定値を用いた突発検知
 - ▶出力誤差の理論値と実験結果の照合