

ストリームデータアルゴリズム

アルゴリズム論

中野

Note 11 ストリームデータアルゴリズム

2020.6.15 作成

2020.7.3 update 7.6 2021.5.31 6.27 6.28

近似カウンタ

https://en.wikipedia.org/wiki/Approximate_counting_algorithm

大量のイベントを少しのメモリでカウントしたい

例 送信元IPアドレスごとのパケットのカウント

IPアドレスは32bit or 8bit x 4



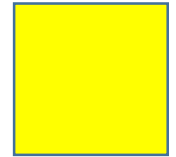
近似カウンタ

https://en.wikipedia.org/wiki/Approximate_counting_algorithm

大量のイベントを少しのメモリでカウントしたい

素朴な方法

確率 $1/1000$ でカウントを増加する



イベント1000万回 \Rightarrow カウンタの値 だいたい1万になる
イベント1500回するとき \Rightarrow カウンタの値 1500? 2? 1.5? 1? 0?
相対誤差 (真の値の割合) 大きいかも

Morrisのアルゴリズム



1,2,4,8,16,...回のいずれであるかを記憶する。

8bitなら 2^{255} まで（確率的に）カウントできる。 $2^{255} > 5.78 \times 10^{76}$

つまり 2^{c-1} の $c=1,2,3,\dots$ を記憶する。ただし $c=0$ のとき0回とする。

カウンタの増加は確率的にする。

例えば、 c が

$0 \Rightarrow 1$ のときは確率1でカウンタを増加する。1回発生

$1 \Rightarrow 2$ のときは確率1でカウンタを増加する。2回発生

$2 \Rightarrow 3$ のときは確率 $1/2$ でカウンタを増加する。確率 $1/2$ でそのまま。4回発生

$3 \Rightarrow 4$ のときは確率 $1/4$ でカウンタを増加する。確率 $3/4$ でそのまま。8回発生

$4 \Rightarrow 5$ のときは確率 $1/8$ でカウンタを増加する。確率 $7/8$ でそのまま。16回発生

カウンタの値が3になったとき

実際のイベントの発生回数の期待値は

x回目	1	2	3	4	5...		
	0	0	0			$3 \times 1/2$	イベント3回の確率は1/2
	0	0	X	0		$4 \times 1/4$	イベント4回の確率は1/4
	0	0	X	X	0	$5 \times 1/8$	
	0	0	X	X	X	$6 \times 1/16$	
						

期待値 $E = 3 \times 1/2 + 4 \times 1/4 + 5 \times 1/8 + 6 \times 1/16 + \dots$

$$E/2 = 3 \times 1/4 + 4 \times 1/8 + 5 \times 1/16 + 6 \times 1/32 + \dots$$

$$E/2 = 3 \times 1/2 + 1 \times 1/4 + 1 \times 1/8 + 1 \times 1/16 + 1 \times 1/32 + \dots$$

$$= 3 \times 1/2 + 1/2$$

$$= 2$$

$$E = 4$$

イベントが2回発生するときカウンタが2（2回目）の確率

$$00 \quad 1 \times 1 = 1$$

イベントが3回発生するときカウンタが2（2回目）の確率

$$00x \quad 1 \times 1 \times 1/2 = 1/2$$

イベントが4回発生するときカウンタが2（2回目）の確率

$$00xx \quad 1 \times 1 \times 1/2 \times 1/2 = 1/4$$

イベントが5回発生するときカウンタが2（2回目）の確率

$$00xxx \quad 1 \times 1 \times 1/2 \times 1/2 \times 1/2 = 1/8$$

。。。。

イベントが10回発生するときカウンタが2（2回目）の確率

$$00xxx \cdots \quad 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \cdots 1/2 = 1/256$$



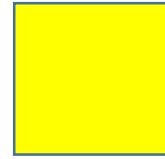
イベントが7回発生するときカウンタが3の確率



7回発生するときカウンタの値が3, つまり $4=2^3-1$ (4回目) になる確率は

$$\begin{aligned} & (00xxxx0) 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \times 1/2 \times 1/2 + \\ & (00xxx0x) 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \times 1/2 \times 3/4 + \\ & (00xx0xx) 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \times 3/4 \times 3/4 + \\ & (00x0xxx) 1 \times 1 \times 1/2 \times 1/2 \times 3/4 \times 3/4 \times 3/4 + \\ & (000xxxx) 1 \times 1 \times 1/2 \times 3/4 \times 3/4 \times 3/4 \times 3/4 \\ & = (16+24+36+54+81)/512 = 211/512 \end{aligned}$$

イベントが7回発生するときカウンタが4の確率



$$\begin{aligned} & (00xxx00) 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \times 1/2 \times 1/4 + \\ & (00xx0x0) 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \times 3/4 \times 1/4 + \\ & (00x0xx0) 1 \times 1 \times 1/2 \times 1/2 \times 3/4 \times 3/4 \times 1/4 + \\ & (000xxx0) 1 \times 1 \times 1/2 \times 3/4 \times 3/4 \times 3/4 \times 1/4 + \\ & (00xx00x) 1 \times 1 \times 1/2 \times 1/2 \times 1/2 \times 1/4 \times 3/4 + \\ & (00x0x0x) 1 \times 1 \times 1/2 \times 1/2 \times 3/4 \times 1/4 \times 7/8 + \\ & (000xx0x) 1 \times 1 \times 1/2 \times 3/4 \times 3/4 \times 1/4 \times 7/8 + \\ & (00x00xx) 1 \times 1 \times 1/2 \times 1/2 \times 1/4 \times 7/8 \times 7/8 + \\ & (000x0xx) 1 \times 1 \times 1/2 \times 3/4 \times 1/4 \times 7/8 \times 7/8 + \\ & (0000xxx) 1 \times 1 \times 1/2 \times 1/4 \times 7/8 \times 7/8 \times 7/8 = ?? \end{aligned}$$

次の問題を考えます

Finding Frequent Items in Data Streams

<http://archive.dimacs.rutgers.edu/Workshops/WGUnifyingTheory/Slides/cormode.pdf>

大量のデータのストリーム。データ n 個。

メモリの制限あり。カウンタ k 個。

その時点までで n/k 回以上出現したデータを収集したい。

Finding Frequent Items in Data Streams

<http://archive.dimacs.rutgers.edu/Workshops/WGUnifyingTheory/Slides/cormode.pdf>

大量のデータのストリーム。

メモリの制限あり。カウンタ k 個。

その時点までで n/k 以上出現したデータを収集したい。

高々 k 個のカウンタを維持する。

新データが計数中のとき \Rightarrow そのカウンタを $+1$ 。

新データが計数中でないとき

(場合1) カウンタの個数 $< k$ \Rightarrow 新カウンタを準備して 1 にする。

(場合2) カウンタの個数 $= k$ \Rightarrow すべてのカウンタを -1

Finding Frequent Items in Data Streams

<http://archive.dimacs.rutgers.edu/Workshops/WGUnifyingTheory/Slides/cormode.pdf>

高々 k 個のカウンタを維持する。

新データが計数中のとき \Rightarrow そのカウンタを $+1$ 。

新データが計数中でないとき

(場合1) カウンタの個数 $< k$ \Rightarrow 新カウンタを準備して 1 にする。

(場合2) カウンタの個数 $= k$ \Rightarrow **すべてのカウンタを -1**

(解析)

-1 発生 \Rightarrow カウンタに記録されたちょうど k 個分のデータと新データの合計 $k+1$ 個のデータを削除

つまり -1 は **高々 $n/(k+1)$ 回**おこる

カウンタ中のデータは高々 $n/(k+1)$ 回しか -1 されない！

よって $(n/(k+1)) + 1 \leq n/k$ **回以上出現したデータは生きのこるはず**

カウントのエラーは高々 $n/(k+1)$

例

$n = 10^{12}$ 1テラ個

$k = 10^6$ 1メガ個

各カウンタ 5Byte=40bit だと

$2^{40} = 1.09 \times 10^{12}$ 回出現までカウントできる。

メモリサイズは $10^6 \times 5 = 5$ メガ Byte。

$10^{12} / 10^6 = 10^6 = 1$ メガ回以上出現のデータを数える。

(0.0001%以上出現するデータを近似的に数える。)

エラーは最大でも $10^{12} / 10^6 = 10^6 = 1$ メガ回 (1テラの $1/10^6$ です)