

MATA54 Estrutura de Dados e Algoritmos II – 2024.2

# Count–Min Sketch

**Alunos**

Gustavo de Oliveira Ferreira

# CONTEXTO

Tenho uma aplicação que está recebendo uma **stream** de eventos e precisamos contar fazer algumas análises

- Quais hashtags estão sendo mais usadas em uma rede social em um certo momento?
- Quais dados devo dar preferência para fazer cache?

# Abordagem Simples

## **Hash Tables?**

Sempre que tiver a ocorrência de um evento posso incrementar o valor dele na tabela

# Abordagem Simples

## **Hash Tables?**

Sempre que tiver a ocorrência de um evento posso incrementar o valor dele na tabela

## **Árvores Binárias de Busca?**

Armazenar as chaves, com os valores de frequência maneira ordenada

# Abordagem Simples

Mas e se estiver analisando uma quantidade massiva de dados?

## **Hash Tables**

Complexidade espacial  $O(n)$

## **Árvores Binárias de Busca?**

Complexidade espacial  $O(n)$

# Abordagem Simples

Mas e se estiver analisando uma quantidade massiva de dados?

## **Hash Tables**

Complexidade espacial  $O(n)$

## **Árvores Binárias de Busca?**

Complexidade espacial  $O(n)$

**Complexidade espacial linear**

# Abordagem Simples

Mas e se estiver analisando uma quantidade massiva de dados?

## **Hash Tables**

Complexidade espacial  $O(n)$

## **Árvores Binárias de Busca?**

Complexidade espacial  $O(n)$

**Mas e se o espaço disponível for restrito?**

**MOTIVAÇÃO**



# Motivação

Analisar stream de eventos **potencialmente** infinitos

Quantidade massiva de dados

Limitação de espaço para guardar tais dados

Análise rápida

# Motivação

Analisar stream de eventos **potencialmente** infinitos

Quantidade massiva de dados

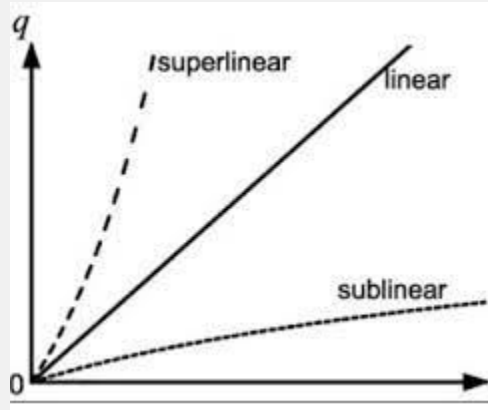
Limitação de espaço para guardar tais dados

Análise rápida

**Precisamos de um estrutura de dados que possua espaço sub-linear.**

# Motivação

Estrutura de dados que possua **espaço sub-linear**.



# COUNT-MIN SKETCH

# Count-min Sketch

Algoritmo Probabilístico para **estimativa** de valores

Pode ser usado para:

- Calcular frequência simples de um evento
- Identificar elementos frequentes
- Computar quantidades

# Count-min Sketch

Realiza operações usando uma matriz de contadores e uma série de funções hash.

hash 1	Mtr 11	Mtr 12	Mtr 13	...	Mtr 1M
hash 1	Mtr 21	Mtr 22	Mtr 23	...	Mtr 2M
...	...	...	...	...	...
hash n	Mtr M1	Mtr M2	Mtr M3	...	Mtr MN

# Tabela CMS

	0	1	2	3	4	5	6
H1	0	0	0	0	0	0	0
H2	0	0	0	0	0	0	0
H3	0	0	0	0	0	0	0
H4	0	0	0	0	0	0	0

# Valores

Stream: [ A, B, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	<b>1</b>	0	0	0	0	0
H2	0	0	0	0	0	0	0
H3	0	0	0	0	0	0	0
H4	0	0	0	0	0	0	0

## Valores

Stream: [ **A**, B, K, A, A, K, S... ]

**H1(A) = 1**, H2(A) = 6, H3(A) = 3, H4(A) = 1

H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6

H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6

H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1



## Tabela CMS

	0	1	2	3	4	5	6
H1	0	1	0	0	0	0	0
H2	0	0	0	0	0	0	<b>1</b>
H3	0	0	0	0	0	0	0
H4	0	0	0	0	0	0	0

## Valores

Stream: [ **A**, B, K, A, A, K, S... ]

$H1(A) = 1$ ,  **$H2(A) = 6$** ,  $H3(A) = 3$ ,  $H4(A) = 1$

$H1(B) = 1$ ,  $H2(B) = 2$ ,  $H3(B) = 4$ ,  $H4(B) = 6$

$H1(K) = 3$ ,  $H2(K) = 4$ ,  $H3(K) = 1$ ,  $H4(K) = 6$

$H1(S) = 6$ ,  $H2(S) = 2$ ,  $H3(S) = 4$ ,  $H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	1	0	0	0	0	0
H2	0	0	0	0	0	0	1
H3	0	0	0	1	0	0	0
H4	0	0	0	0	0	0	0

## Valores

Stream: [ **A**, B, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, \mathbf{H3(A) = 3}, H4(A) = 1$

$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	1	0	0	0	0	0
H2	0	0	0	0	0	0	1
H3	0	0	0	1	0	0	0
H4	0	1	0	0	0	0	0

## Valores

Stream: [ **A**, B, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, \mathbf{H4(A) = 1}$

$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	<b>2</b>	0	0	0	0	0
H2	0	0	0	0	0	0	1
H3	0	0	0	1	0	0	0
H4	0	1	0	0	0	0	0

## Valores

Stream: [ A, **B**, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

**$H1(B) = 1$** ,  $H2(B) = 2, H3(B) = 4, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	2	0	0	0	0	0
H2	0	0	<b>1</b>	0	0	0	1
H3	0	0	0	1	0	0	0
H4	0	1	0	0	0	0	0

## Valores

Stream: [ A, **B**, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, \mathbf{H2(B) = 2}, H3(B) = 4, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	2	0	0	0	0	0
H2	0	0	1	0	0	0	1
H3	0	0	0	1	<b>1</b>	0	0
H4	0	1	0	0	0	0	0

## Valores

Stream: [ A, **B**, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, H2(B) = 2, \mathbf{H3(B) = 4}, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	2	0	0	0	0	0
H2	0	0	1	0	0	0	1
H3	0	0	0	1	1	0	0
H4	0	1	0	0	0	0	<b>1</b>

## Valores

Stream: [ A, **B**, K, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, H2(B) = 2, H3(B) = 4, \mathbf{H4(B) = 6}$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	2	0	<b>1</b>	0	0	0
H2	0	0	1	0	<b>1</b>	0	1
H3	0	<b>1</b>	0	1	1	0	0
H4	0	1	0	0	0	0	<b>2</b>

## Valores

Stream: [ A, B, **K**, A, A, K, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, H2(A) = 2, H3(A) = 4, H4(A) = 6$

**$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$**

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$



## Tabela CMS

	0	1	2	3	4	5	6
H1	0	<b>3</b>	0	1	0	0	0
H2	0	0	1	0	1	0	<b>2</b>
H3	0	1	0	<b>2</b>	1	0	0
H4	0	<b>2</b>	0	0	0	0	2

## Valores

Stream: [ A, B, K, **A**, A, K, S... ]

**H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1**

H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6

H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6

H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1

# Tabela CMS

	0	1	2	3	4	5	6
H1	0	<b>4</b>	0	1	0	0	0
H2	0	0	1	0	1	0	<b>3</b>
H3	0	1	0	<b>3</b>	1	0	0
H4	0	<b>3</b>	0	0	0	0	2

# Valores

Stream: [ A, B, K, A, **A**, K, S... ]

**H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1**

H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6

H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6

H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	4	0	<b>2</b>	0	0	0
H2	0	0	1	0	<b>2</b>	0	3
H3	0	<b>2</b>	0	3	1	0	0
H4	0	3	0	0	0	0	<b>3</b>

## Valores

Stream: [ A, B, K, A, A, **K**, S... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$

**$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$**

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	4	0	2	0	0	<b>1</b>
H2	0	0	<b>2</b>	0	2	0	3
H3	0	2	0	3	<b>2</b>	0	0
H4	0	<b>4</b>	0	0	0	0	3

## Valores

Stream: [ A, B, K, A, A, K, **s...** ]

$$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$$

$$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$$

$$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$$

$$\mathbf{H1(s) = 6, H2(s) = 2, H3(s) = 4, H4(s) = 1}$$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	<b>4</b>	0	2	0	0	1
H2	0	0	2	0	2	0	<b>3</b>
H3	0	2	0	<b>3</b>	2	0	0
H4	0	<b>4</b>	0	0	0	0	3

Quantas vezes **A** foi recebido no stream?

$$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$$

$$\min [4, \mathbf{3}, \mathbf{3}, 4] \Rightarrow 3 \text{ vezes}$$

## Valores

Stream: [ A, B, K, A, A, K, S... ]

$$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$$

$$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$$

$$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$$

$$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$$

## Tabela CMS

	0	1	2	3	4	5	6
H1	0	4	0	2	0	0	<b>1</b>
H2	0	0	<b>2</b>	0	2	0	3
H3	0	2	0	3	<b>2</b>	0	0
H4	0	<b>4</b>	0	0	0	0	3

Quantas vezes **S** foi recebido no stream?

$$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$$

$$\min [1, 2, 2, 4] \Rightarrow 1 \text{ vez}$$

## Valores

Stream: [ A, B, K, A, A, K, S... ]

$$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$$

$$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$$

$$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$$

$$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$$

## Tabela CMS - Colisão

	0	1	2	3	4	5	6
H1	0	4	0	2	0	0	<b>2</b>
H2	<b>1</b>	0	2	0	2	0	3
H3	<b>1</b>	2	0	3	2	0	0
H4	0	4	<b>1</b>	0	0	0	3

Mas e se nós adicionarmos um novo valor a stream, que a função hash incremente a posição do menor valor de outro elemento?

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

$\min [2, 2, 2, 4] \Rightarrow$  No máximo 2 vezes

## Adicionando E ao stream

Stream: [ A, B, K, A, A, K, S, **E**... ]

$H1(A) = 1, H2(A) = 6, H3(A) = 3, H4(A) = 1$

$H1(B) = 1, H2(B) = 2, H3(B) = 4, H4(B) = 6$

$H1(K) = 3, H2(K) = 4, H3(K) = 1, H4(K) = 6$

$H1(S) = 6, H2(S) = 2, H3(S) = 4, H4(S) = 1$

**$H1(E) = 6, H2(E) = 0, H3(E) = 0, H4(E) = 2$**

# Count-min Sketch Pseudocódigo

```
1: procedure CountMinSketch(depth, width)
2:   function add(item)
3:     while counter < depth do
4:       buckets  $\leftarrow$  getHashBuckets(item, depth, width)
5:       M tr[counter][buckets[counter]]  $\leftarrow$  +1
6:       counter  $\leftarrow$  +1
7:     end while
8:   end function
```



# Count-min Sketch Pseudocódigo

```
1:  procedure CountMinSketch(depth, width)
9:      function estimateCount(item)
10:           $r \leftarrow \text{MAXVALUE}$ 
11:          while counter < depth do
12:               $\text{val} \leftarrow M \text{ tr}[\text{counter}][\text{buckets}[\text{counter}]]$ 
13:               $r \leftarrow \text{Min}(r, \text{val})$ 
14:              counter  $\leftarrow +1$ 
15:          end while
16:      return r
17:  end function
```

**TEORIA**

# Probabilidade

O CMS é uma estrutura com parâmetros  $(\varepsilon, \delta)$  representada por uma matriz bi-dimensional.

$\varepsilon$  representa o fator de erro na estimativa de frequência (precisão da estrutura de dados)

$\delta$  representa a probabilidade de erro (a chance de que a estimativa exceda o erro especificado por  $\varepsilon$ )

# Probabilidade

A matriz é formada por duas listas com largura **w** e profundidade **d**.

Dados os parâmetros  $(\epsilon, \delta)$ , temos que:

# Probabilidade

A matriz é formada por duas listas com largura **w** e profundidade **d**.

Dados os parâmetros  $(\epsilon, \delta)$ , temos que:

$$\text{largura } w = \lceil e/\epsilon \rceil \qquad \text{profundidade } w = \lceil \ln 1/\delta \rceil$$

# Probabilidade

Dado:

- **m** elementos adicionados
- **$\hat{a}$**  estimativa de ocorrência de um elemento  $x$
- **a** quantidade real de ocorrência de  $x$

Temos que:

$$\Pr(\hat{a} \leq a + \varepsilon * m) \geq 1 - \delta$$

# Probabilidade

Dado:

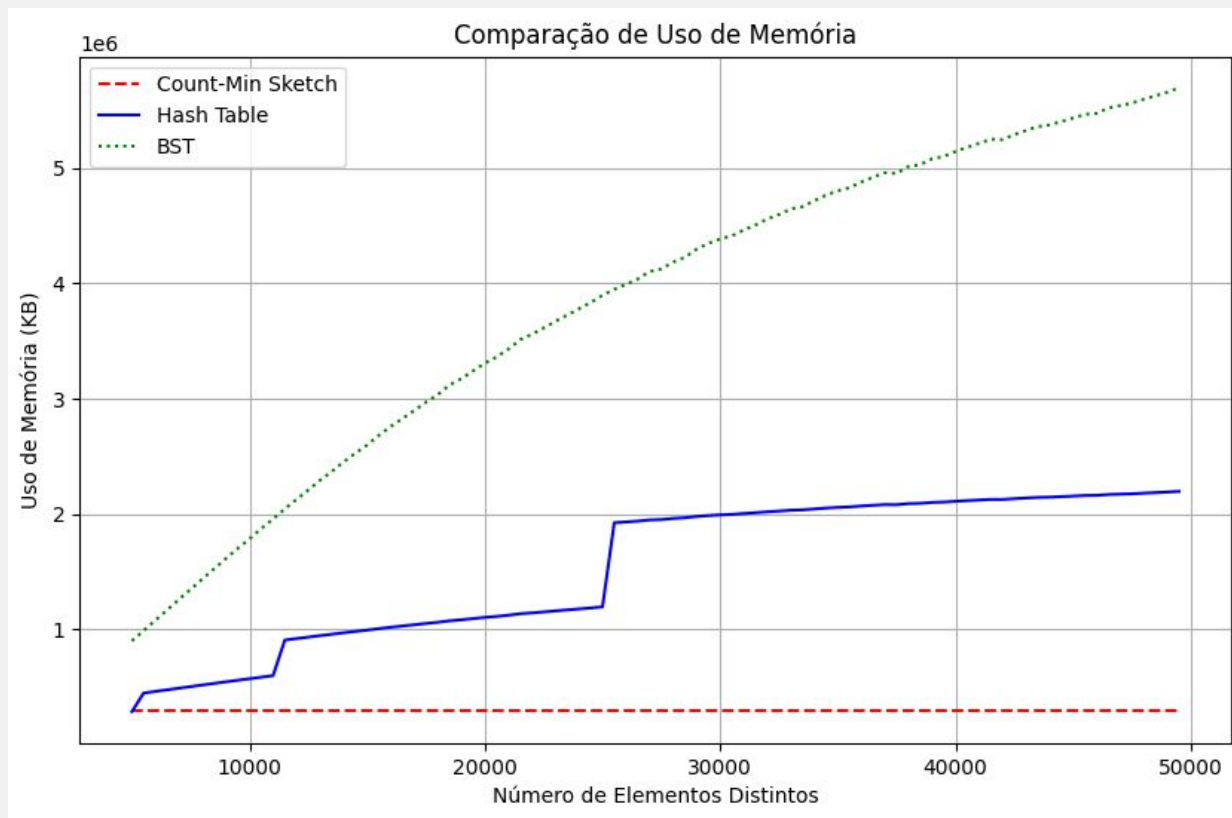
- **m** elementos adicionados
- **$\hat{a}$**  estimativa de ocorrência de um elemento x
- **a** quantidade real de ocorrência de x

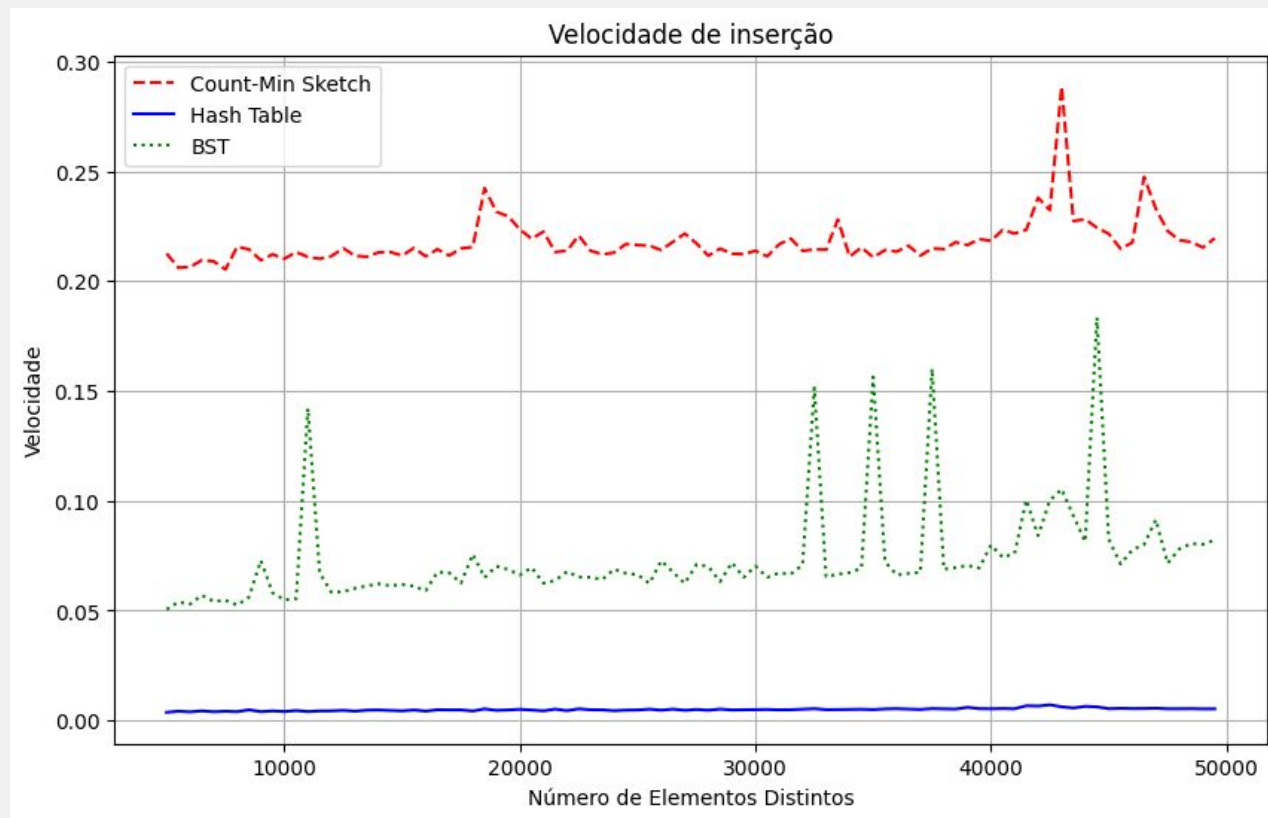
Temos que:

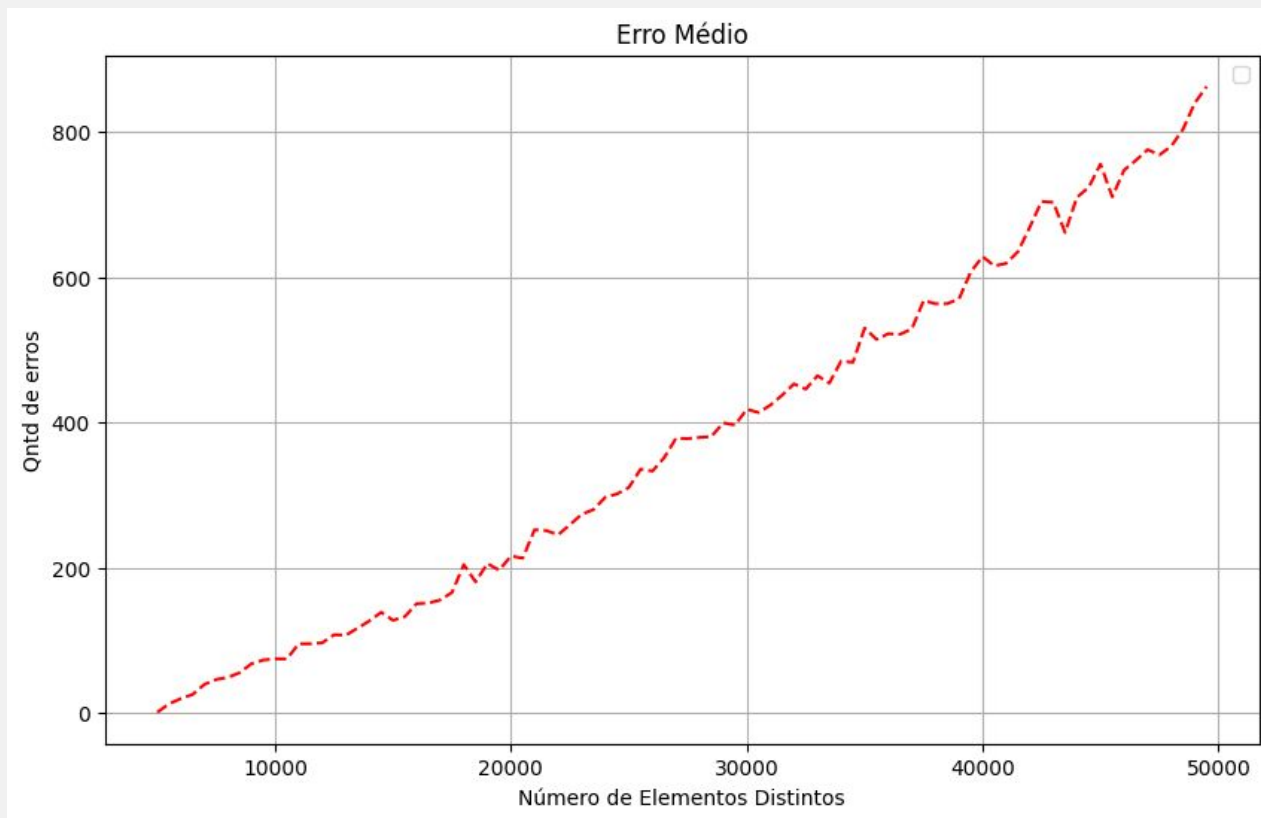
$$\Pr(\hat{a} \leq a + \varepsilon * m) \geq 1 - \delta$$

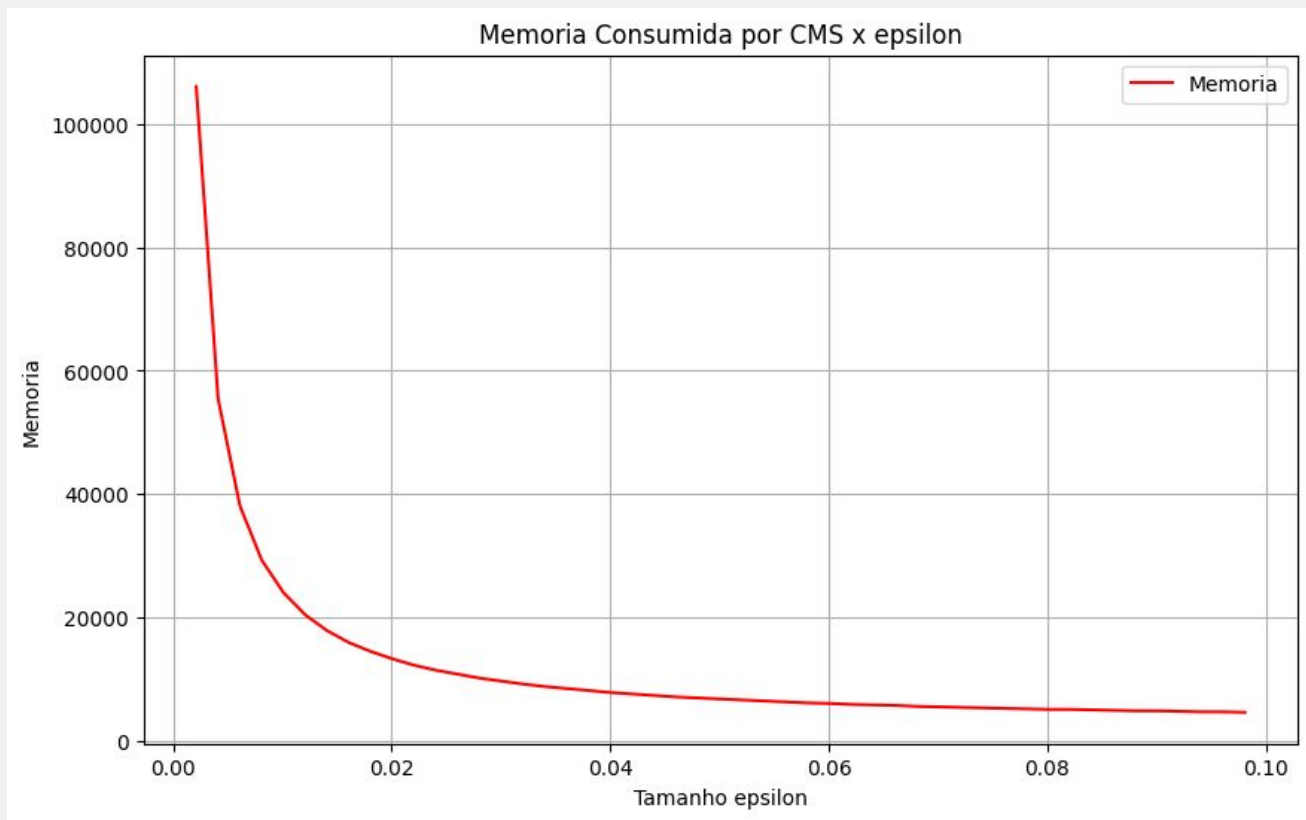
# ANÁLISE PRÁTICA











# Referências

CORMODE, G.; MUTHUKRISHNAN, S. **An improved data stream summary: the count-min sketch and its applications.** Journal of Algorithms, v. 55, n. 1, p. 58–75, abr. 2005.

CORMODE, G.; MUTHUKRISHNAN, S. **What's hot and what's not: tracking most frequent items dynamically.** ACM Transactions on Database Systems, v. 30, n. 1, p. 249–278, mar. 2005.

SCHWARZ, K. **Count-Min Sketches.** Disponível em: <<https://web.stanford.edu/class/archive/cs/cs166/cs166.1206/lectures/10/Slides10.pdf>>. Acesso em: 20 jan. 2025.

**Obrigado!**