

Algoritmo Count-Min Sketch

Introdução e História

O **Count-Min Sketch (CMS)** é uma técnica de resumo probabilístico (ou estrutura de dados de resumo) que permite contar frequências aproximadas de elementos em fluxos de dados massivos, consumindo muito menos memória do que armazenar todas as ocorrências.

O algoritmo foi proposto por **Graham Cormode** e **S. Muthukrishnan** em 2005, no artigo seminal:

Cormode, G., & Muthukrishnan, S. (2005). An improved data stream summary: The count-min sketch and its applications. *Journal of Algorithms*, 55(1), 58-75.

Esse artigo consolidou o CMS como uma solução prática para o problema de contar frequências de itens em fluxos de dados (data streams), onde o volume é tão grande que não cabe na memória principal para armazenamento exato.

Por que o Count-Min Sketch surgiu?

Antes do CMS, o problema de contagem de frequências em streams era tratado por estruturas como **hash tables** ou **arrays** tradicionais, que demandam espaço proporcional ao número de elementos distintos (chamados de *universe size*). Em cenários de grandes volumes, como monitoramento de tráfego de rede, logs de acessos web ou análise de dados em tempo real, manter um contador exato para cada elemento era impraticável.

Técnicas anteriores tentavam equilibrar espaço e precisão, mas enfrentavam limitações:

- **Hash tables exatas:** Precisam de muito espaço.
- **Reservoir sampling e outras técnicas de amostragem:** Podem não estimar frequências com precisão.
- **Bloom filters:** Respondem perguntas de pertencimento, não frequências.

CMS aparece como um compromisso: estimar a frequência com erros controlados, usando espaço sublinear, e com garantia probabilística.

Funcionamento do Count-Min Sketch

O CMS é uma matriz 2D de contadores com dimensões $d \times w$, onde:

- d é o número de funções hash independentes,
- w é o tamanho da "largura" da tabela (quantidade de colunas).

Cada função hash mapeia um elemento para uma das w colunas.

Inserção: Para cada elemento que chega no stream, incrementa-se um contador em cada uma das d linhas, na coluna correspondente ao hash daquele elemento.

Consulta: Para estimar a frequência do elemento, calcula-se o valor mínimo entre todos os d contadores hashados para aquele elemento.

Essa escolha do mínimo — que dá o nome "count-min" — evita superestimação causada por colisões de hash.

Propriedades matemáticas

- O erro estimado na contagem é sempre positivo (não subestima).
- A probabilidade do erro ultrapassar um limite é controlada por d e w .
- Especificamente, dado erro ϵ e probabilidade δ , podemos configurar $w = \lceil \epsilon / \delta \rceil$ e $d = \lceil \ln(1/\delta) / \epsilon \rceil$.

Vantagens do Count-Min Sketch

- **Eficiência espacial:** Requer espaço muito menor que armazenar os contadores exatos.
- **Velocidade:** Operações de inserção e consulta são rápidas — tempo constante $O(d)$, normalmente pequeno.
- **Simplicidade:** Fácil de implementar e adaptar a várias aplicações.
- **Garantias teóricas:** Erro e probabilidade controlados e previsíveis.
- **Aplicações diversas:** Análise de tráfego, bancos de dados, compressão, monitoramento, aprendizado de máquina, etc.

Cormode e Muthukrishnan destacam que:

"The count-min sketch is a compact summary structure for streaming data that can be used in various tasks such as point queries, inner product estimation, and more."

Desvantagens e limitações

- **Erro sempre positivo:** O CMS só superestima, nunca subestima, o que pode ser problemático para algumas aplicações.
- **Colisões de hash:** Podem levar a erros maiores se as funções hash não forem suficientemente independentes.
- **Não armazena itens explicitamente:** Não é possível recuperar os itens mais frequentes sem técnicas auxiliares.
- **Parâmetros sensíveis:** A escolha de d e w impacta diretamente precisão e memória.
- **Não captura frequências exatas:** Apenas aproximações, o que pode não ser aceitável para aplicações críticas.

O que era feito antes do CMS?

Antes da chegada do CMS, o problema de contagem em streams era abordado principalmente por:

- **Métodos exatos**, como tabelas de hash tradicionais, que não escalavam para dados grandes.

- **Sketches probabilísticos anteriores**, como o **Count Sketch** (Charikar et al., 2002), que estimam frequências, mas com distribuições de erro diferentes.
- **Amostragem** e outras técnicas estatísticas que não garantiam eficiência ou precisão rigorosa para todo o universo de elementos.

O CMS trouxe um equilíbrio entre simplicidade, eficiência e garantias teóricas, o que o tornou uma das ferramentas fundamentais em processamento de dados massivos.

Funcionamento detalhado do algoritmo Count-Min Sketch

Estrutura de dados

O Count-Min Sketch é representado por uma matriz C de inteiros com d linhas e w colunas, onde:

- d é o número de funções hash independentes.
- w é o tamanho da tabela para cada função hash.

Cada célula da matriz armazena um contador que representa a frequência aproximada de elementos que foram mapeados para aquela posição.

Inicialização

- Inicializamos todos os contadores da matriz C com zero.
- Selecionamos d funções hash independentes h_1, h_2, \dots, h_d , cada uma mapeando um item x para um índice entre 0 e $w-1$.

Essas funções hash são essenciais para distribuir uniformemente os elementos na matriz, reduzindo colisões.

Inserção de um elemento

Quando um elemento x aparece no fluxo de dados (stream), o algoritmo executa:

Para cada linha $i \in \{1, 2, \dots, d\}$:

1. Calcula o índice da coluna: $j = h_i(x)$.
2. Incrementa o contador $C[i, j]$ em 1 (ou na quantidade associada, caso o elemento venha com um peso ou frequência).

Ou seja:

$$\text{Para } i = 1 \text{ a } d: \quad C[i, h_i(x)] \leftarrow C[i, h_i(x)] + 1$$

Assim, para cada elemento, incrementamos d contadores, um por linha, na coluna definida pela função hash daquela linha.

Consulta (estimativa da frequência)

Para estimar a frequência do elemento x , o algoritmo faz:

Para cada linha $i \in \{1, 2, \dots, d\}$:

1. Calcula o índice da coluna: $j = h_i(x)$.
2. Recupera o valor do contador $C[i, j]$.

O valor estimado para a frequência de x , \hat{f}_x , é o **mínimo** desses valores:

\$\$

$$\hat{f}_x = \min_{1 \leq i \leq d} C[i, h_i(x)]$$

\$\$

Por que o mínimo? Porque cada contador $C[i, h_i(x)]$ pode conter incrementos causados por colisões (outros elementos mapeados para a mesma posição), e portanto pode superestimar a frequência real. A tomada do mínimo das linhas ajuda a aproximar melhor a frequência verdadeira, minimizando o impacto das colisões.

Análise de erro

O erro na estimativa é um ponto fundamental do CMS.

Suponha:

- f_x é a frequência verdadeira do elemento x .
- \hat{f}_x é a frequência estimada pelo CMS.

Então:

\$\$

$$f_x \leq \hat{f}_x \leq f_x + \frac{\epsilon}{N}$$

\$\$

onde:

- N é o número total de elementos inseridos no sketch (tamanho do stream).
- ϵ é o erro máximo permitido (fração do total N).

A probabilidade de que a estimativa ultrapasse esse limite é:

\$\$

$$P(\hat{f}_x > f_x + \frac{\epsilon}{N}) \leq \delta$$

\$\$

Os parâmetros ϵ (erro) e δ (falha de probabilidade) são configurados definindo:

- $w = \lceil \frac{e}{\epsilon} \rceil$ (largura da matriz, onde e é o número de Euler)
- $d = \lceil \ln(1 / \delta) \rceil$ (profundidade da matriz)

Essas configurações garantem que o erro seja pequeno com alta probabilidade.

Complexidade

- **Espaço:** $O(d \times w)$, que é sublinear em relação ao tamanho do universo ou número total de elementos distintos.
- **Tempo de inserção:** $O(d)$, porque precisa atualizar d posições.
- **Tempo de consulta:** $O(d)$, pois verifica d posições para estimar a frequência.

Na prática, d e w são pequenos (tipicamente $d = 5$ a 10) para oferecer bons trade-offs entre precisão e uso de memória.

Exemplo prático simplificado

Imagine que queremos monitorar a frequência de palavras em um stream enorme de texto.

- Escolhemos $d = 3$ funções hash e $w = 1000$ colunas.
- Para cada palavra que chega, calculamos 3 índices hash e incrementamos os contadores nas posições correspondentes.
- Para consultar a frequência da palavra "chatGPT", consultamos os 3 contadores e pegamos o menor valor.

Mesmo que algumas outras palavras causem colisões e incrementem esses contadores, a tomada do mínimo reduz o impacto da superestimação.

Resumo do algoritmo

Passo	Descrição
Inicialização	Criar matriz $d \times w$ e definir funções hash
Inserção	Incrementar d contadores, um por linha hashada
Consulta	Estimar frequência pegando o mínimo dos d contadores hashados

Claro! Vamos aprofundar bastante no uso prático do **Count-Min Sketch (CMS)**, explicando suas aplicações reais no dia a dia, em empresas e projetos, com exemplos concretos de setores que se beneficiam dessa estrutura.

Aplicações do Count-Min Sketch no Mundo Real

O Count-Min Sketch é uma estrutura de dados fundamental em cenários onde há necessidade de processar grandes volumes de dados em tempo real, de forma eficiente em espaço e tempo, e com tolerância a erros controlados. Seu uso se espalhou amplamente em empresas de tecnologia, telecomunicações, sistemas distribuídos e análise de dados massivos.

1. Monitoramento e análise de tráfego de rede

Contexto

Operadoras de internet, provedores de serviços em nuvem e grandes data centers precisam monitorar o tráfego de rede em tempo real para identificar padrões, detectar ataques, e otimizar o uso dos recursos.

Aplicação do CMS

- **Contagem de pacotes ou fluxos:** É preciso contabilizar rapidamente a frequência de IPs, protocolos, portas ou fluxos específicos que passam pela rede.
- **Deteção de ataques DDoS:** Ao monitorar os IPs com alta frequência de requisições, o CMS pode ajudar a identificar picos suspeitos.
- **Balanceamento de carga e QoS:** Conhecer as frequências de acessos a servidores ou serviços auxilia no balanceamento automático.

Exemplos

- **Cisco** e outras empresas de equipamentos de rede utilizam técnicas de sketching, como o CMS, para manter estatísticas de tráfego em alta velocidade, onde não seria possível guardar todas as informações em memória.
- Projetos como o **NetFlow** da Cisco empregam resumos compactos para contabilizar fluxos e identificar os mais frequentes.

2. Sistemas de recomendação e análise de logs

Contexto

Plataformas de comércio eletrônico, redes sociais e serviços de streaming coletam bilhões de eventos de usuários, como cliques, visualizações e interações.

Aplicação do CMS

- **Contar eventos de usuários:** Frequência de cliques em produtos, vídeos assistidos, termos de busca.
- **Identificação rápida de tendências:** Detectar palavras-chave, hashtags ou itens que estão se tornando populares.
- **Compressão de logs:** Reduzir o espaço necessário para armazenar dados estatísticos históricos.

Exemplos

- **Twitter** usa técnicas semelhantes a CMS para manter contagem de hashtags e termos populares em tempo real, permitindo identificar tendências emergentes.

- **Netflix e Spotify** aplicam sketches para analisar padrões de uso e recomendações, sem precisar armazenar contagem exata para cada item.
-

3. Bancos de dados e sistemas distribuídos

Contexto

Sistemas de bancos de dados NoSQL, caches distribuídos e sistemas de busca precisam gerenciar grandes volumes de chaves e consultas, muitas vezes em ambientes distribuídos e sob restrições de memória.

Aplicação do CMS

- **Estimativa rápida de frequência de chaves:** Para otimizar caching, indexação ou balanceamento de carga.
- **Sincronização eficiente:** Reduzir dados trocados entre nós de um cluster ao enviar resumos compactos.
- **Controle de acesso e rate limiting:** Estimar quantas requisições um usuário fez para aplicar limites.

Exemplos

- O **Apache Cassandra**, banco NoSQL amplamente utilizado, pode integrar estruturas de resumo para monitorar consultas e otimizar armazenamento.
 - Em sistemas como **Google Bigtable** e **Amazon DynamoDB**, técnicas semelhantes ajudam a identificar "hot keys" (chaves acessadas com alta frequência) para balancear cargas.
-

4. Segurança da Informação

Contexto

Soluções de segurança precisam analisar grandes volumes de eventos de logs para detectar padrões suspeitos, anomalias e ataques.

Aplicação do CMS

- **Deteção de IPs maliciosos ou domínios suspeitos:** Contagem rápida das frequências para identificar fontes de ataque.
- **Análise de logs:** Agregar dados de acesso para detectar padrões incomuns.

Exemplos

- Ferramentas de SIEM (Security Information and Event Management) usam sketches para sumarizar grandes volumes de logs e detectar ameaças em tempo real.
 - Projetos acadêmicos e comerciais na área de segurança aplicam CMS para melhorar a performance de detecção de ataques DDoS e varreduras de portas.
-

5. Publicidade e análise de cliques (Ad Tech)

Contexto

Empresas de publicidade digital monitoram bilhões de cliques, impressões e eventos para otimizar campanhas e detectar fraudes.

Aplicação do CMS

- **Detecção de fraude em cliques:** Monitorar IPs e usuários para identificar cliques suspeitos repetidos.
- **Contagem aproximada de impressões:** Para medir o alcance de anúncios em tempo real.
- **Otimização de campanhas:** Identificar quais anúncios têm melhor desempenho.

Exemplos

- Plataformas como **Google Ads** e **Facebook Ads** integram técnicas de sketching para analisar grandes volumes de dados em tempo real.
- Empresas de Ad Tech menores também usam CMS para garantir que a análise de dados seja rápida e consuma pouca memória.

6. Machine Learning e Big Data

Contexto

Modelos de aprendizado de máquina que processam dados em larga escala precisam de técnicas para pré-processamento e sumarização de dados.

Aplicação do CMS

- **Extração de características:** Contar frequências aproximadas de tokens ou atributos em grandes datasets.
- **Redução de dimensionalidade:** Usar contadores compactos para gerar features que representem o dado original.
- **Aprendizado em streaming:** Manter estatísticas atualizadas em tempo real.

Exemplos

- Frameworks de aprendizado de máquina distribuído, como o **Apache Spark**, podem integrar técnicas de sketch para sumarizar dados distribuídos.
- Aplicações em NLP (Processamento de Linguagem Natural) usam CMS para contagem rápida de n-gramas em grandes coleções de texto.

Resumo dos Benefícios para Empresas

Benefícios	Exemplos de Setores
------------	---------------------

Benefícios	Exemplos de Setores
Uso reduzido de memória	Data centers, provedores de nuvem
Velocidade no processamento	Redes de telecom, Ad Tech
Estimativas confiáveis com erro controlado	Bancos de dados, segurança
Aplicação em tempo real	Plataformas sociais, sistemas IoT

Casos famosos e menções na literatura

- **Cormode e Muthukrishnan (2005)** discutem aplicações em monitoramento de rede, sistemas de bancos de dados e análise de logs.
- A **Twitter** divulgou em apresentações técnicas o uso de sketches para identificar tendências em hashtags, dada a escala de bilhões de tweets diários.
- Projetos do **Google** e **Facebook** integram variações do CMS em suas infraestruturas para monitoramento e análise em tempo real.

O Count-Min Sketch é hoje uma peça-chave em infraestruturas que lidam com big data e streaming. Sua capacidade de fornecer contagens aproximadas com espaço reduzido permite que empresas escalem operações que seriam impossíveis com abordagens tradicionais.

A aplicação do CMS permite detectar tendências, otimizar recursos, detectar fraudes e monitorar sistemas complexos em tempo real, em setores que vão desde redes até publicidade, passando por segurança e machine learning.

A eficiência e flexibilidade dessa estrutura fazem dela uma ferramenta essencial no arsenal do cientista de dados e engenheiro de software moderno.

Claro! Vou fornecer implementações básicas do **Count-Min Sketch** tanto em **C** quanto em **Python**. Ambas são simples, ilustrativas e focadas em conceitos, para você entender o funcionamento e poder adaptar para aplicações maiores.

Implementação do Count-Min Sketch

1. Implementação em C

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>

#define MAX_D 5          // número de funções hash (profundidade)
#define MAX_W 1000       // largura da matriz (tamanho de cada linha)

// Estrutura do Count-Min Sketch
typedef struct {
```

```

    int d;                // profundidade
    int w;                // largura
    unsigned int **table; // matriz de contadores
} CountMinSketch;

// Função hash simples baseada em djb2 (com seed para independência)
unsigned int hash(unsigned int seed, const char *str) {
    unsigned int hash = 5381 + seed;
    int c;
    while ((c = *str++))
        hash = ((hash << 5) + hash) + c; /* hash * 33 + c */
    return hash;
}

// Inicializa o CMS
CountMinSketch* cms_create(int d, int w) {
    CountMinSketch *cms =
    (CountMinSketch*)malloc(sizeof(CountMinSketch));
    cms->d = d;
    cms->w = w;
    cms->table = (unsigned int **)malloc(d * sizeof(unsigned int*));
    for (int i = 0; i < d; i++) {
        cms->table[i] = (unsigned int *)calloc(w, sizeof(unsigned int));
    }
    return cms;
}

// Incrementa a frequência do elemento
void cms_update(CountMinSketch *cms, const char *element) {
    for (int i = 0; i < cms->d; i++) {
        unsigned int idx = hash(i, element) % cms->w;
        cms->table[i][idx]++;
    }
}

// Estima a frequência do elemento
unsigned int cms_query(CountMinSketch *cms, const char *element) {
    unsigned int min = ~0; // maior valor possível de unsigned int
    for (int i = 0; i < cms->d; i++) {
        unsigned int idx = hash(i, element) % cms->w;
        if (cms->table[i][idx] < min)
            min = cms->table[i][idx];
    }
    return min;
}

// Libera a memória
void cms_free(CountMinSketch *cms) {
    for (int i = 0; i < cms->d; i++)
        free(cms->table[i]);
    free(cms->table);
    free(cms);
}

```

```
// Exemplo de uso
int main() {
    CountMinSketch *cms = cms_create(MAX_D, MAX_W);

    cms_update(cms, "apple");
    cms_update(cms, "banana");
    cms_update(cms, "apple");
    cms_update(cms, "orange");
    cms_update(cms, "banana");
    cms_update(cms, "apple");

    printf("Frequencia estimada de apple: %u\n", cms_query(cms,
"apple")); // Esperado ~3
    printf("Frequencia estimada de banana: %u\n", cms_query(cms,
"banana")); // Esperado ~2
    printf("Frequencia estimada de orange: %u\n", cms_query(cms,
"orange")); // Esperado ~1
    printf("Frequencia estimada de grape: %u\n", cms_query(cms,
"grape")); // Esperado 0 (ou muito próximo)

    cms_free(cms);
    return 0;
}
```

Explicação rápida:

- Usamos 5 funções hash diferentes, simuladas pelo seed da função hash.
- A matriz \$d \times w\$ é alocada dinamicamente.
- Atualizamos o CMS incrementando o contador para cada hash.
- Na consulta, pegamos o mínimo entre as 5 posições.
- Para strings, usamos a função hash djb2 com diferentes seeds para simular múltiplas funções hash.

2. Implementação em Python

```
class CountMinSketch:
    def __init__(self, d=5, w=1000):
        import hashlib
        self.d = d
        self.w = w
        self.table = [[0]*w for _ in range(d)]
        self.salts = [str(i).encode('utf-8') for i in range(d)] # para
funções hash diferentes

    def _hash(self, x, i):
        import hashlib
        h = hashlib.md5(self.salts[i] + x.encode('utf-8')).hexdigest()
        return int(h, 16) % self.w
```

```

def update(self, x):
    for i in range(self.d):
        idx = self._hash(x, i)
        self.table[i][idx] += 1

def query(self, x):
    min_count = float('inf')
    for i in range(self.d):
        idx = self._hash(x, i)
        if self.table[i][idx] < min_count:
            min_count = self.table[i][idx]
    return min_count

# Exemplo de uso
cms = CountMinSketch()

cms.update("apple")
cms.update("banana")
cms.update("apple")
cms.update("orange")
cms.update("banana")
cms.update("apple")

print("Frequencia estimada de apple:", cms.query("apple"))    # Esperado
~3
print("Frequencia estimada de banana:", cms.query("banana"))  # Esperado
~2
print("Frequencia estimada de orange:", cms.query("orange"))  # Esperado
~1
print("Frequencia estimada de grape:", cms.query("grape"))    # Esperado
0 (ou próximo)

```

Explicação rápida:

- Usamos `hashlib.md5` com diferentes "salts" para simular múltiplas funções hash.
- A matriz é uma lista de listas.
- O processo de inserção e consulta segue o padrão do CMS.
- Fácil de entender e usar em protótipos ou sistemas com menor necessidade de desempenho.

Se quiser, posso ajudar a expandir essas implementações para incluir:

- Suporte a pesos (incrementos maiores que 1)
- Serialização para salvar o sketch em arquivo
- Uso em streaming real
- Integração com outras estruturas de dados

Conclusão

O **Count-Min Sketch** é uma estrutura de dados fundamental no campo de algoritmos para *data streams*. Criado para resolver o problema prático de contar frequências em fluxos muito grandes, ele fornece uma aproximação eficiente com garantias probabilísticas de erro.

Sua importância se traduz em inúmeras aplicações práticas que envolvem monitoramento em tempo real, análise de grandes volumes e sistemas que não comportam armazenar todos os dados explicitamente.

Como ressalta Cormode e Muthukrishnan (2005):

"The count-min sketch is a versatile tool enabling approximate queries over massive data streams with strong performance guarantees."

Contudo, deve ser aplicado com consciência de suas limitações, especialmente em contextos que exigem contagens exatas ou evitam erros positivos.

Referências principais

- Cormode, G., & Muthukrishnan, S. (2005). An improved data stream summary: The count-min sketch and its applications. *Journal of Algorithms*, 55(1), 58–75.
- Charikar, M., Chen, K., & Farach-Colton, M. (2002). Finding frequent items in data streams. *International Colloquium on Automata, Languages, and Programming (ICALP)*.
- Mitzenmacher, M., & Upfal, E. (2005). *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge University Press.