Aula: Estruturas de Dados – Algoritmo Trie

Introdução

A estrutura de dados **Trie** é uma das mais importantes e eficientes para armazenamento e recuperação de conjuntos de strings, especialmente quando o problema envolve buscas por prefixos ou palavras completas em grandes volumes de dados textuais. Também conhecida como **árvore prefixo**, a Trie possui aplicações essenciais em sistemas de busca, correção ortográfica, auto-completar, dicionários digitais e compressão de dados.

1. História e Motivação

A Trie foi introduzida formalmente pelo cientista informático René de la Briandais em 1959, em seu artigo "File Searching Using Variable Length Keys" (Communications of the ACM, 1959). O termo "Trie" vem da palavra "retrieval", apesar de muitos associarem a palavra erroneamente a "tree" (árvore).

Naquela época, sistemas computacionais já precisavam armazenar grandes volumes de dados textuais, como índices de livros, listas telefônicas e bancos de dados. Antes da Trie, estruturas comuns para esse fim eram:

- Árvores binárias de busca (BST): eficientes para buscas em geral, mas não especializadas para strings, especialmente para operações de prefixo.
- **Listas encadeadas ou arrays ordenados**: simples, porém muito ineficientes para grandes volumes e buscas frequentes.
- Tabelas hash: rápidas para buscas exatas, mas não eficientes para buscas por prefixos, além de não manterem ordenação.

Segundo Knuth (The Art of Computer Programming, Volume 3, 1997), as Trie são particularmente eficientes para "o problema de recuperação de palavras em dicionários ou índices, especialmente quando se busca por prefixos".

2. O que era feito antes da Trie?

PROFESSEUR: M.DA ROS

Antes do advento da Trie, o armazenamento e a busca de strings em grandes bases de dados textuais dependiam principalmente de:

- Árvores de busca binárias (BSTs): onde as palavras eram inseridas e buscadas em ordem lexicográfica. No entanto, o custo da busca dependia do comprimento da palavra, e a estrutura não explorava o compartilhamento de prefixos entre as palavras.
- Hashing direto: para buscas exatas, mas que não suportava buscas por prefixos nem ordenação.
- **Vetores ou listas ordenadas**: facilitavam buscas binárias, mas o custo de inserção e busca era elevado para grandes datasets.

Essas limitações motivaram a busca por uma estrutura que aproveitasse o fato de que muitas palavras compartilham prefixos, economizando espaço e tempo.

3. Como funciona a Trie?

A Trie é uma árvore em que:

- Cada nó representa um caractere (ou símbolo) da string.
- O caminho da raiz até um nó representa um prefixo da palavra armazenada.
- Palavras são armazenadas nas folhas ou em nós marcados como finais.
- Os filhos de cada nó representam possíveis extensões do prefixo.

Assim, para buscar uma palavra, o algoritmo percorre a Trie caractere por caractere, seguindo o caminho correspondente. Se a palavra existir, a busca termina em um nó final; caso contrário, a busca falha rapidamente.

4. Vantagens da Trie

- Busca eficiente por prefixo: permite encontrar todas as palavras que começam com um dado prefixo em tempo proporcional ao comprimento do prefixo, independente do número total de palavras na estrutura.
- Economia de espaço para dados com prefixos compartilhados: ao reutilizar nós comuns a vários prefixos, evita redundância.
- Operações rápidas de inserção e remoção: também em tempo proporcional ao comprimento da palavra.
- Suporte para ordenação lexicográfica automática: devido à organização em árvore, a travessia pré-ordem da Trie produz palavras ordenadas lexicograficamente.

Autores como Cormen et al. (Introduction to Algorithms, 3rd Edition, 2009) destacam que a Trie "é uma escolha natural para problemas de recuperação rápida de cadeias de caracteres e manipulação de conjuntos de strings".

5. Desvantagens da Trie

- Alto consumo de memória: especialmente se o alfabeto for grande, já que cada nó pode ter até | Σ| filhos, onde Σ é o alfabeto. Isso pode causar overhead de espaço, sobretudo para dicionários pequenos.
- Implementação mais complexa: comparada a estruturas lineares, exige maior cuidado com gerenciamento de nós e referências.
- Não é eficiente para chaves muito curtas: onde outras estruturas, como tabelas hash, podem ser mais rápidas.
- Problemas com dados dispersos: em casos onde poucas palavras compartilham prefixos, o ganho de espaço é pequeno.

6. Aplicações práticas e variações

Além da Trie padrão, várias variações foram desenvolvidas para otimizar espaço e desempenho, como:

- Compressed Trie (Radix Tree): que compacta cadeias de caracteres comuns em um único nó.
- Suffix Trie: usado para indexação de substrings em análise de texto e bioinformática.
- Ternary Search Trie: combina a Trie com árvores de busca ternárias para economizar espaço.

Onde o algoritmo Trie é usado no dia a dia?

1. Sistemas de busca e autocompletar:

- Quando você digita algo no Google, no seu navegador ou em aplicativos de mensagem, a sugestão automática de palavras ou frases é frequentemente baseada em estruturas como Trie.
- A Trie permite encontrar rapidamente todas as palavras que começam com o prefixo digitado, entregando sugestões instantâneas.

2. Dicionários e corretor ortográfico:

 Aplicativos de edição de texto e corretores ortográficos usam Trie para armazenar grandes dicionários de palavras. Isso permite verificar rapidamente se uma palavra existe e sugerir correções com base em prefixos semelhantes.

3. Sistemas de roteamento de redes e IP:

Em redes, Tries especializadas (como Patricia tries) são usadas para armazenar rotas IP,
 buscando rapidamente qual é o caminho correto para um endereço IP com base em
 prefixos.

4. Compressão de dados e processamento de texto:

- Algoritmos de compressão, como o LZ78, usam variações de Trie para identificar padrões repetidos em textos.
- Em bioinformática, Tries são usadas para indexar sequências de DNA e RNA, facilitando buscas rápidas em grandes bases genéticas.

Como podemos melhorar o algoritmo Trie para uso prático?

Embora a Trie seja eficiente, ela pode ser aprimorada para enfrentar problemas reais de uso, principalmente relacionados ao consumo de memória e velocidade:

1. Compressão de nós – Radix Trees / Compressed Trie

- Em vez de armazenar um caractere por nó, compactamos sequências de caracteres que não têm bifurcação em um único nó.
- Isso reduz drasticamente o número de nós e o consumo de memória.
- Por exemplo, em vez de ter um caminho com nós 'c' → 'a' → 't', podemos ter um único nó com a palavra "cat".

2. Uso de Trie Ternária (Ternary Search Trie)

- Uma Trie ternária combina características de árvores binárias e Tries, reduzindo o número de ponteiros por nó para 3 (menor, igual, maior).
- Isso reduz o uso de memória em alfabetos grandes, tornando o algoritmo mais escalável.

3. Implementação com estruturas mais eficientes de armazenamento

- Em vez de usar arrays para armazenar filhos (que podem ser grandes e esparsos), usamos tabelas hash ou mapas para armazenar apenas os filhos existentes.
- Isso reduz o espaço desperdiçado com ponteiros nulos.

4. Lazy deletion e balanceamento

- Implementar estratégias para remoção eficiente de palavras sem reconstruir a Trie inteira.
- Balancear a Trie para reduzir a profundidade e melhorar a velocidade de busca.

5. Persistência e estruturas híbridas

- Para aplicações em bancos de dados ou sistemas distribuídos, Tries podem ser persistidas em disco usando estruturas híbridas que combinam Tries com B-Trees, para garantir performance mesmo em dados massivos.
- No dia a dia, Tries aparecem em buscas rápidas por palavras e prefixos, correção ortográfica, redes, compressão e bioinformática.
- Para melhorar, compressão de nós, uso de Tries ternárias, armazenamento eficiente e técnicas de balanceamento são essenciais para deixar a Trie escalável e prática.

Implementação da Trie em Python

```
class TrieNode:
   def __init__(self):
        self.children = {} # dicionário para filhos: char -> TrieNode
        self.is_end_of_word = False # marca o fim de uma palavra
class Trie:
   def init (self):
        self.root = TrieNode()
   def insert(self, word: str):
        node = self.root
        for char in word:
            if char not in node.children:
                node.children[char] = TrieNode()
            node = node.children[char]
        node is end of word = True
   def search(self, word: str) -> bool:
        node = self.root
```

```
for char in word:
    if char not in node.children:
        return False
    node = node.children[char]
    return node.is_end_of_word

# Exemplo de uso:
    trie = Trie()
    trie.insert("carro")
    trie.insert("casa")
    print(trie.search("carro")) # True
    print(trie.search("casa")) # False
    print(trie.search("casa")) # True
```

Implementação da Trie em C

```
#include <stdio.h>
#include <stdlib.h>
#include <stdbool.h>
#define ALPHABET_SIZE 26
// Estrutura do nó da Trie
typedef struct TrieNode {
    struct TrieNode *children[ALPHABET SIZE];
    bool is_end_of_word;
} TrieNode;
// Função para criar um novo nó da Trie
TrieNode* createNode() {
    TrieNode *node = (TrieNode *)malloc(sizeof(TrieNode));
    node->is end of word = false;
    for (int i = 0; i < ALPHABET SIZE; i++) {
        node->children[i] = NULL;
    return node;
}
// Função para inserir uma palavra na Trie
void insert(TrieNode *root, const char *word) {
    TrieNode *node = root;
    for (int i = 0; word[i] != '\setminus 0'; i++) {
        int index = word[i] - 'a'; // assume somente letras minúsculas
        if (node->children[index] == NULL) {
            node->children[index] = createNode();
        node = node->children[index];
    }
```

```
node->is_end_of_word = true;
}
// Função para buscar uma palavra na Trie
bool search(TrieNode *root, const char *word) {
    TrieNode *node = root;
    for (int i = 0; word[i] != '\setminus 0'; i++) {
        int index = word[i] - 'a';
        if (node->children[index] == NULL) {
            return false;
        node = node->children[index];
    return node->is_end_of_word;
}
// Exemplo de uso
int main() {
    TrieNode *trie = createNode();
    insert(trie, "carro");
    insert(trie, "casa");
    printf("%d\n", search(trie, "carro")); // 1 (true)
    printf("%d\n", search(trie, "cas")); // 0 (false)
    printf("%d\n", search(trie, "casa")); // 1 (true)
    // Note: Não implementamos liberação da memória aqui para
simplicidade.
    return 0;
}
```

Observações

- Ambas as implementações assumem o alfabeto de letras minúsculas a-z.
- Em Python usamos um dicionário para filhos, o que é mais flexível.
- Em C usamos um vetor fixo de tamanho 26, o que limita o alfabeto, mas é eficiente.
- Em aplicações reais, é importante implementar a liberação da memória (free em C) e melhorias para lidar com outros caracteres (acentos, letras maiúsculas, etc.).

7. Conclusão

PROFESSEUR: M.DA ROS

A estrutura Trie representa um avanço significativo para o problema clássico da recuperação eficiente de cadeias de caracteres, especialmente quando a operação de busca por prefixos é fundamental. Desde sua proposição por René de la Briandais, vem sendo estudada, otimizada e aplicada em áreas que vão desde editores de texto até sistemas de busca e bioinformática.

Apesar das limitações no consumo de memória, seu desempenho na manipulação de grandes volumes de dados textuais com prefixos comuns justifica seu uso em várias aplicações. Conforme Knuth (1997), "a Trie é uma das estruturas de dados fundamentais para algoritmos de recuperação de informação".

Referências

- De la Briandais, R. (1959). *File Searching Using Variable Length Keys*. Communications of the ACM, 3(5), 295-298.
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (2009). *Introduction to Algorithms* (3rd ed.). MIT Press.
- Knuth, D. E. (1997). *The Art of Computer Programming, Volume 3: Sorting and Searching* (2nd ed.). Addison-Wesley.
- Fredkin, E. (1960). Trie Memory. Communications of the ACM, 3(9), 490-499.