

TP1: Banco de Dados com Arvores B+

Giovani Amaral Barcelos

5 de maio de 2016

1 Introdução

O objetivo desse trabalho prático é construir um banco de dados simplificado usando uma árvore B+, que conterá operações simples de busca de chave, inserção de registro e impressão por nível. Além disso, temos o detalhe que entre duas operações somente a raiz ficará em memória, ao passo que todos os outros nós ficarão em disco.

A solução apresentada consiste na implementação de uma árvore B+ e, para manter os nós em disco, foi usado serialização e desserialização dos mesmos, sempre que for preciso carregar algum nó para a memória ou persisti-lo em disco.

2 Solução do problema

2.1 Primeira parte - árvore B+

O primeiro grande desafio do problema é implementar uma árvore B+. Alguns detalhes permitiram fazer isso com mais facilidade: 1) as chaves dos registros sempre são únicas e 2) sempre são valores numéricos. Desta forma, foi adaptado o algoritmo da árvore B do [1] para se tornar uma árvore B+, atendendo nossa necessidade. Caso o leitor não saiba, a diferença é que, na árvore B+, só temos registros nas folhas, e cada folha tem um apontador para seu irmão a direita, facilitando caminhamento sequencial. No código da árvore, modificamos algumas coisas para atender ao critério de somente a raiz em memória.

2.1.1 Fluxo externo - leitura de dados

Para lidar com as operações básicas, temos um loop que lê todas as linhas do arquivo de entradas, e para cada linha é identificada a operação (*add*, *search* ou *dump*) e então é chamada a função correspondente.

Algorithm 1: Pseudo-código do fluxo principal

```
1 for linha ∈ arquivo do
2   | campos ← processar_linha(linha, numero_campos, operacao)
3   | if operacao == "add" then
4   |   | registro ← cria_registro(campos)
5   |   | inserir(raiz, arquivo_arvore, registro)
6   | end
7   | if operacao == "search" then
8   |   | registro ← busca(arquivo_arvore, raiz, chave, numero_campos)
9   |   | Print(registro)
10  | end
11  | if operacao == "dump" then
12  |   | busca_largura(raiz, arquivo_saida, arquivo_arvore, numeroCampos)
13  | end
14 end
```

2.1.2 Inserção

O código da inserção pode ser aproximado pelo código abaixo. Pode-se notar que sua base se parece com o código do do [1], mas aqui nós desalocamos um nó, evitando que tenha mais de um nó em memória. Essa alteração é feita na linha 11, que faz com que sempre que precisamos dividir um nó, nós desaloquemos o que não vamos usar mais nessa inserção. Caso formos usar novamente, isso será decidido dentro da função de inserir com espaço.

Algorithm 2: Pseudo-código da função inserir

Data: nodo* raiz, FILE* arquivo, Registro* reg
Result: nodo* raiz

```
1 // Salvamos a ordem em uma variável para facilitar o uso
2 ordem ← no → ordem
3 r ← raiz
4 // Máximo de registros
5 if r → numeroRegistros == ordem then
6   nodo * s ← aloca_nodo()
7   raiz ← s
8   s → numeroRegistros ← 0
9   separa_filho(arquivo, s, 1, r, reg → numeroCampos)
10  // Diferente do Cormem, precisamos desalocar devido ao nosso critério
    especial
11  desaloco_nodo(r, reg → numeroCampos)
12  // Passamos o nodo no qual vamos tentar inserir e 1 para indicar que é a
    raiz
13  insere_nodo_com_espaco(s, arquivo, reg, 1)
14 else
15   insere_nodo_com_espaco(r, arquivo, reg, 1)
16 end
17 return raiz
```

No algoritmo abaixo, nós estaremos inserindo um registro em um nó que tem espaço. Se esse nó for uma folha, inserimos nele. Se não, vamos para seu filho que contenha o valor correto para inserção. A grande alteração desse algoritmo que faz com que ele seja compatível com nosso requerimento, é que sempre, ao terminar uma inserção, nós desalocamos um nó caso ele não seja uma raiz(linha 18). Caso precisemos caminhar na árvore para inserir mais embaixo, iremos desalocando os nós a medida que formos descendo. Isso é visto na linha 39 e 45.

Algorithm 3: Pseudo-codigo da funcao `inserir_nodo_com_espaco`

Data: `nodo*` raiz, `FILE*` arquivo, `Registro*` reg, `int` is_raiz

Result: `void`

```
1 // Salvamos a ordem em uma variável para facilitar o uso
2 ordem ← no → ordem
3 i ← r → numeroRegistros
4 if r → folha then
5     // Movemos os registros para o lado
6     while i ≥ 1 and reg → x → chaves[i] do
7         | x → chaves[i + 1] ← x → chaves[i]
8         | x → registros[i + 1] ← x → registros[i]
9         | i ← i − 1
10    end
11    // Insere registro na posição correta
12    x → chaves[i + 1] ← reg → chave
13    x → registros[i + 1] ← reg
14    x → numeroRegistros ← x → numeroRegistros + 1
15    serializar(fp, x, x → registros[i + 1] → numeroCampos, x → offset)
16    // Desalocamos o nó, caso ele não seja uma raiz
17    if !is_raiz then
18        | desaloca_nodo(x, reg → numeroCampos)
19    end
20 end
21 else
22     // Decidimos para qual filho vai este registro
23     while i ≥ 1 and reg → x → chaves[i] do
24         | i ← i − 1
25     end
26     i++
27     nodo* no
28     aloca_memoria(no)
29     // Recuperamos o filho correto do arquivo
30     deserializar(arquivo, x → offsetFilhos[i], no)
31     // Se estiver cheio, vamos dividi-lo
32     if no → numeroRegistros == ordem then
33         | separa_filho(arquivo, x, i, no, reg → numeroCampos)
34         | // Checamos se vamos inserir no da esquerda ou direita
35         | if reg → chave > x → chaves[i] then
36             | i++
37         | end
38     end
39     desaloca_nodo(no, reg → numeroCampos)
40     // Recuperamos novamente o nó, caso o nó destino tenha mudado após a
41     // divisão
42     aloca_memoria(no)
43     deserializar(arquivo, x → offsetFilhos[i], no)
44     // Se não é raiz, desalocamos
45     if !is_raiz then
46         | desaloca_nodo(x, reg → numeroCampos)
47     end
48     // Agora passamos 0, pois o no sendo passado não é a raiz
49     insere_nodo_com_espaco(no, arquivo, reg, 0)
50 end
```

Como último algoritmo importante da parte de inserção, temos o `separa_filho`, que recebe um nó e um pai, e separa o filho em dois.

Algorithm 4: Pseudo-codigo da funcao `separa_filho`

```

Data: nodo* x, FILE* arquivo, nodo* y, int i, int numero_campos
Result: void
1  Seja ordem a ordem da árvore
2  aloca_memoria(z)
3  if  $y \rightarrow folha$  then
4    | separa_folha( $y, z, ordem$ )
5  end
6  else
7    | separa_interno( $y, z, ordem$ )
8  end
9  // Movemos os valores dos offsets do filhos para a direita, para inserir o
   novo valor do novo nó gerado na divisão
10 for  $j \leftarrow x \rightarrow numeroRegistros + 1$  to  $j \geq 1 + 1$  do
11   |  $x \rightarrow offsetFilhos[j + 1] \leftarrow x \rightarrow offsetFilhos[j]$ 
12 end
13  $x \rightarrow offsetFilhos[i + 1] \leftarrow z \rightarrow offset$ 
14 // Movemos agora as chaves
15 for  $j \leftarrow x \rightarrow numeroRegistros$  to  $j \geq 1$  do
16   |  $x \rightarrow chaves[j + 1] \leftarrow x \rightarrow chaves[j]$ 
17 end
18  $x \rightarrow chaves[i] \leftarrow y \rightarrow chaves[(ordem/2) + 1]$ 
19  $x \rightarrow numeroRegistros ++$ 
20 serializar(arquivo,  $x$ , numeroCampos,  $x \rightarrow offset$ )
21 serializar(arquivo,  $y$ , numeroCampos,  $y \rightarrow offset$ )
22 serializar(arquivo,  $z$ , numeroCampos,  $z \rightarrow offset$ )
23 desaloca_nodo( $z$ , numeroCampos)

```

2.1.3 Busca em largura - Impressão

Os algoritmos acima nos ajudarão a analisar a complexidade geral do nosso algoritmo e especifica em funções. Vamos ver aqui apenas mais 4 funções importantes: 2 da árvore (busca em largura e pesquisa) e 2 das operações em disco (serialização e desserialização).

Algorithm 5: Pseudo-codigo da funcao `busca em largura`

```

Data: nodo* raiz, FILE* arquivo_saida, FILE* arquivo_arvore, int numero_campos
Result: void
1  Seja f uma fila encadeada
2  Enfileira(raiz)
3  while  $f \rightarrow vazia$  do
4    |  $n \leftarrow Desenfileira()$ 
5    | for  $filho \in n$  do
6      | if  $n \rightarrow folha$  then
7        | | Enfileira(filho)
8      | end
9    | end
10   | Print( $n \rightarrow registros$ )
11   | desaloca_nodo( $n$ )
12 end

```

2.1.4 Busca de chave

Nossa próxima e última operação da árvore será a busca:

Algorithm 6: Pseudo-codigo da funcao busca de chave

Data: nodo* raiz, int k, FILE* arquivo_arvore, int numero_campos
Result: Registro*

```
1 Seja registro_encontrado um ponteiro para registro
2 while  $i \leq x \rightarrow \text{numeroRegistros} \text{ AND } k \Rightarrow x \rightarrow \text{chaves}[i]$  do
3   | i++
4 end
5 if  $x \rightarrow \text{folha}$  then
6   | if  $k == x \rightarrow \text{chaves}[i - 1]$  then
7     |  $\text{registro\_encontrado} \leftarrow x \rightarrow \text{registros}[i - 1]$ 
8     | desaloca_nodo(x, numeroCampos)
9     | return registro_encontrado
10  | end
11  | else
12    | desaloca_nodo(x, numeroCampos)
13    | return NULL
14  | end
15 end
16 aloca(no)
17 desserializar(arquivo_arvore,  $x \rightarrow \text{offsetFilhos}[i]$ , no)
18 desaloca_nodo(x, numeroCampos)
19 return busca(arquivo_arvore, no, k, numeroCampos);
```

2.2 Segunda parte - Serialização e Desserialização

Nosso grande desafio foi manter apenas a raiz em memória entre uma operação e outra. Só é permitido ter a raiz em memória, e para isso, não poderíamos ter apontadores apontando para os filhos, mas sim guardamos offsets que indicariam a posição do filho no arquivo. Desta forma, podemos recuperar um nó sabendo a sua posição no arquivo.

2.2.1 Serialização

Na hora de serializar, era necessário guardar o maior espaço que um nodo poderia ocupar. Como exemplo, vamos imaginar uma folha com 3 registro em uma árvore ordem 5. Se salvamos apenas 3 registros no disco e esta árvore recebe mais um, precisaríamos escrever além do espaço reservado para ela no arquivo, ou seja, sobreescreveríamos informações importantes para outros nodos.

Desta forma, a melhor solução foi preencher com 0's onde poderia vir a existir um registro.

O código da serialização pode ser tangenciado pelo seguinte pseudo código:

Algorithm 7: *Serializacao*

Data: nodo* raiz, long long unsigned offset, FILE* arquivo_arvore, int numero_campos

Result: void

```
1 Seja metadata um vetor de longlongunsigned com 5 posições
2 metadata[0] ← no → numeroRegistros
3 metadata[1] ← no → folha
4 metadata[2] ← no → ordem
5 metadata[3] ← no → offset
6 metadata[4] ← no → prox
7 EscreveArquivo(arquivo_arvore, metadata)
8 // Serializamos agora o vetor de chaves
9 serializa_vetor_chaves(arquivo_arvore, no)
10 // Serializamos agora o vetor de offsets dos filhos
11 serializa_vetor_offsetFilhos(arquivo_arvore, no)
12 // Se o nó é uma folha, serializamos também seus registros
13 if no → folha then
14 |   serializa_matriz_registros(arquivo_arvore, no → registros, no → ordem, no →
15 |   folha, numeroCampos, no → numeroRegistros)
15 end
```

Você pode estar se perguntando o motivo de não colocarmos, por exemplo, o vetor de chaves no metadata. Isso se dá pois ele é um ponteiro, e não queremos serializar um endereço, mas sim uma lista de valores. Por isso precisamos serializar um a um.

2.3 Desserializacao

Por fim, dado um offset, essa função deve ser capaz de ir até essa posição no arquivo e carregar o nó que está armazenado lá, através da leitura e conversão de seus bytes.

Algorithm 8: *Desserializacao*

Data: nodo* no, long long unsigned offset, FILE* arquivo_arvore

Result: void

```
1 Colocamos o ponteiro do arquivo na posição indicada por offset
2 Seja no um no com mempria previamente alocada
3 no → numeroRegistros ← le_arquivo(longlongunsigned, 1, arquivo_arvore)
4 no → folha ← le_arquivo(longlongunsigned, 1, arquivo_arvore)
5 no → ordem ← le_arquivo(longlongunsigned, 1, arquivo_arvore)
6 no → offset ← le_arquivo(longlongunsigned, 1, arquivo_arvore)
7 no → prox ← le_arquivo(longlongunsigned, 1, arquivo_arvore)
8 aloca_memoria(no → chaves)
9 aloca_memoria(no → offsetFilhos)
10 desserializa_vetor_chaves(arquivo_arvore, no → chaves, no → ordem)
11 desserializa_vetor_offsetFilhos(arquivo_arvore, no → offsetFilhos, no → ordem)
12 // Se for uma folha, lemos os registros
13 if no → folha then
14 |   desserializa_registros(arquivo_arvore, no, no → ordem)
15 end
```

3 Análise teórica do custo assintótico de tempo

Nas análises seguintes, consideraremos M a ordem da árvore. Além disso, chamaremos de Z o tamanho do novo nó criado na separação e Y o tamanho do nó que foi dividido. Desta forma, X é o nó pai dos dois citadas anteriormente. Consideramos também C sendo o número de campos que um registro possui e N sendo o número de registros da árvore.

É importante ressaltar que Y é sempre $M/2$ e Z é sempre $M/2$ ou $(M/2)-1$, mas podemos desconsiderar esse -1 e utilizaremos Y e Z para facilitar o entendimento.

desaloca_nodo - Decidi analisar esta função pois ela é utilizado ao longo de todo o código. Entretanto, sua complexidade é simples, pois ela possui um for que roda o numero de registro, e para cada vez que roda chama `desaloca_registro`, que roda o número de campos C . Logo, nossa complexidade é $O(MC)$.

aloca_nodo - A justificativa desta análise é igual a da anterior. Sua complexidade é $O(MC)$, pois ela chama a função de serializar.

serializar - *Algorithm7* - Nossa função tem algumas operações $O(1)$. Na linha 9, chama uma função que serializa as chaves. Como ela serializa o máximo de chaves que for possível(M), sua complexidade é $O(M)$. Na linha 11 temos a mesma ideia da serialização de chaves, mas temos um offset a mais que chaves, logo roda $M+1$ vezes. No pior caso, se for um registro, ele serializa o número máximo de registros(M)(linha 14) e cada registro possui C campos, logo temos um for que roda C vezes.

No final, a complexidade de nossa função é : $1 + M + M+1 + (M*C) = 2M + MC = M(2+C) = MC$. Complexidade $O(MC)$

Desserialização *Algorithm8* - Essa complexidade se assemelha a complexidade de serialização. Serializamos as chaves ($O(M)$) e offsets ($O(M+1)$). No pior caso de ser um registro, serializamos no máximo M registro e cada um com C campos. Logo, nossa complexidade é de $M + M + 1 + MC = 2M + MC = O(MC)$

separa_interno - temos um for que roda $Z+1$ vezes e o outro que roda Z vezes, totalizando $2Z+1$. Logo a complexidade da nossa função é $O(Z)$;

separa_folha - temos um for que roda Z vezes, logo a complexidade é $O(Z)$

separa_filho - *Algorithm4* - Essa função chama um das duas funções citadas acima(linha 4 ou 7), com complexidade $O(Z)$. Na linha 10, temos um for que desloca os offsets do nó pai que roda o número de registros do nó pai - i . Esse número é, no pior caso, $M - 1$, pois o máximo que o pai pode ter de registro é M e o mínimo que i pode ser é 1. $O(M-1)$

Já na linha 15, temos outro for que roda o mesmo tanto que o for anterior. $O(M-1)$.

Nas linhas 20, 21 e 22 chamamos a função de serialização, com complexidade $O(MC)$.

Não podemos esquecer que é feita uma chamada para `aloca_nodo` e `desaloca_nodo`, e suas complexidades são, respectivamente, $O(MC)$ e $O(MC)$. Logo, a complexidade final é $(M-1)*2 + 3*(MC) + Z + MC + MC$, ou seja, $M - 2 + 5MC + Z = M + MC + Z$ e como Z é $M/2$ e podemos ignorar a constante,temos $2M + MC$ que é igual da ordem **$O(M + MC)$ ou $O(MC)$**

Inserir_nodo_com_espaco *Algorithm5* -

Caso o nó que iremos inserir é uma folha, temos na linha 6 um while que roda, nó maximo, o número de registros em um nó, ou seja, M .

Na linha 15 fazemos uma chamada a serializar, que é $O(MC)$. Na linha 18, desalocamos um nó, que tem complexidade $O(MC)$.

Logo, a complexidade desta parte é $M + MC + MC = M(1 + 2C) = O(2MC) = O(MC)$.

Já na segunda parte, caso o nó que estamos olhando não seja uma folha, possuímos na linha 23 um while que roda no máximo M vezes, que é o número máximo de registros.

Na linha 30 e 42, fazemos uma chamada para desserialização, que tem o custo de MC cada, totalizando $2MC$.

Podemos ter, na linha 33, uma chamada para separa filho, com complexidade $O(MC)$. Na linha 39 e 45 desalocamos um nodo, então temos $2MC$ de custo.

Desta forma, a complexidade da segunda parte é $M + 2MC + 2MC = M(1 + 4C) = 4MC$ ou $O(MC)$.

No pior caso, a altura da nossa arvore e $\log_M(N)$. Entao a segunda parte vai rodar $\log_M(N)$ vezes ate

chegar na folha e então rodar a primeira parte uma vez. Logo nossa complexidade é $(\log_M(n) - 1 * MC) + MC$, ou seja, ordem de complexidade: $O(\log_M(n) * MC)$

Inserir - *Algorithm2*- No pior caso, alocamos um nodo (linha 6), desalocamos um nodo (linha 11), chamamos separa_filhos (linha 9) e chamamos insere_nodo com espaço (linha 13), $MC + MC + MC + \log_M(n) * MC$, que equivale a dizer $O(\log_M(n) * MC)$

Busca por chave - *Algorithm6* Temos um while na linha 2 que roda no máximo M vezes. Caso seja uma folha, temos apenas a complexidade de desalocar um nodo, que é $O(MC)$. Sendo uma folha, nossa complexidade é $M + MC = MC$.

Não sendo uma folha, vamos desserializar(MC) na linha 17 e desalocar(MC) na linha 18, com complexidade $2MC$ ou $O(MC)$.

Vamos rodar a segunda parte $\log_m(n) - 1$ vezes, pois é o número de nós que passaremos até chegar em uma folha. Então nossa complexidade será $(\log_m(n) - 1 * MC) + MC$. Conclusão: $O(\log_m(n) * MC)$

Busca em largura - *Algorithm5* - Temos um while na linha 3 que visita todos os nós da árvore, então sua complexidade é $O(N)$, sendo N o número de nós da árvore.. Na linha 11, deslocamos um registro, complexidade $O(M)$. Todas as operações de fila são $O(1)$. Então a complexidade é $O(NM)$.

3.1 Análise teórica do custo assintótico de espaço

Para fazer essa análise, vamos considerar as principais estruturas NODO e REGISTRO. A estrutura registro possui um long(4 bytes), um unsigned long long(8 bytes) e uma matriz de char de tamanho $30 * C$. Logo, sua complexidade de espaço é

$$4 + 8 + 30C = 30C = O(C)$$

Já a estrutura nodo possui 3 inteiros (12 bytes), um long long unsigned int(8 bytes), um apontador de long long unsigned com tamanho M ($8 * M$), um apontador de long long unsigned com tamanho $M + 1(8 * (M + 1))$ e uma matriz de registros contendo M registros, onde cada um tem complexidade $O(C)$. Se for uma folha, nossa complexidade deverá considerar os registros. Se for um nó interno, ele será menor.

Desta forma, nossa complexidade de espaço total de um nó no pior caso é

$$12 + 8 + 8M + 8M + 8 + MC = 16M + MC = O(MC)$$

Para obter a complexidade total do programa, basta a complexidade obtida acima pelo número de nós ou páginas existentes.

4 Análise de experimentos

Para realizar a análise experimental, foi implementado um gerador de testes que gera instruções baseadas em um número de campos e um campo índice pre determinado. Para medir o tempo de execução do código, foi utilizada o comando time do Ubuntu. Cada teste foi rodado 5 vezes e foi feita uma média com os tempos. Para realizar os testes foi utilizada uma máquina com Linux Ubuntu 15, processador I3 e 6GB de RAM.

A segunda coisa a ser analisada foi a quantidade de bytes resultantes no arquivo que guarda a árvore com suas informações.

A seguinte tabela é fruto de testes com uma árvore com ordem 4 e 4 campos:

Número de registros	Tempo execução	Tamanho do arquivo da árvore	Bytes Alocados
644	0.037	533.224	1,632,554
2015	0.103	1.374.304	5,661,010
2016	0.106	1.374.408	5,664,290
2017	0.149	1.374.408	5,667,066
2018	0.150	1.376.408	5,670,210
2019	0.152	1.376.496	5,673,530
10080	0.594	6.858.600	32,474,738
20160	1.229	13.704.680	69,308,042
17348	1.176	11.796.096	59,037,386

Dobrando a ordem, ficamos com 4 campos e ordem 8:

Número de registros	Tempo execução	Tamanho do arquivo da árvore	Bytes Alocados
644	0.055	517.736	2,301,738
2015	0.098	1.612.120	7,803,650
2016	0.101	1.612.120	7,807,562
2017	0.136	1.612.120	7,811,842
10080	0.605	8.029.528	32,474,738
20160	1.193	16.067.832	89,368,570

Dobrando os campos, ficamos com ordem 4 e 8 campos:

Número de registros	Tempo execução	Tamanho do arquivo da árvore	Bytes Alocados
644	0.041	835.112	2,359,770
2015	0.136	2.599.184	7,958,018
2016	0.167	2.599.288	7,962,514
10080	0.615	1.302.3664	43,713,938
20160	1.366	2.605.1856	91,179,530

Através das análises experimentais, foi possível concluir que, quando precisamos dividir um nó e gravamos somente mais um nó no arquivo, seu tamanho aumenta em 2088 bytes. Além disso, vemos que a complexidade é praticamente linear, pois dobrando o número de registros, os outros valores superam um pouco o dobro dos valores do original

Entretanto, ao dobrar a ordem ou os campos, nossos resultados não dobram. Aumentam, mas não dobram.

5 Bibliografia

Referências

- [1] Thomas H Cormen. *Introduction to algorithms*. MIT press, 2009.