

Aula 10: Busca, Ordenação, Hash

Professor(a): Virgínia Fernandes Mota

<http://www.dcc.ufmg.br/~virginiaferm>

ALGORITMOS E ESTRUTURAS DE DADOS - SETOR DE INFORMÁTICA



Busca

- Vamos discutir diferentes estratégias para efetuar a busca de um elemento em um determinado conjunto de dados.
- Vetores x Árvores Binárias de Busca.

- Dado um vetor *vet* com n elementos, desejamos saber de um determinado elemento *elem* está ou não presente no vetor.
- Forma mais simples: Busca linear

```
1 int busca (int n, int *vet, int elem){
2
3     int i;
4     for (i = 0; i < n; i++){
5         if (elem == vet[i])
6             return i;
7     }
8     return -1;
9 }
```

- Busca linear: Muito ineficiente!
- Complexidade $O(n)$: No pior caso, precisaremos realizar n comparações.
- É possível melhorar?

- Assumindo que os elementos do vetor estão ordenados (ordem crescente).

```
1  int busca_ord(int n, int *vet, int elem){
2      int i;
3
4      for (i = 0; i < n; i++){
5          if (elem == vet[i])
6              return i;
7          else if (elem < vet[i])
8              return -1;
9      }
10     return -1;
11 }
```

- Mas o algoritmo continua sendo linear - $O(n)$.
- No entanto, se os elementos estão ordenados, existe uma busca muito mais eficiente.

- Com os elementos ordenados, podemos aplicar um algoritmo mais eficiente: Busca Binária.

```
1 int busca_bin(int n, int *vet, int elem){
2     int ini = 0;
3     int fim = n-1;
4     int meio;
5
6     while (ini <= fim){
7         meio = (ini + fim) / 2;
8         if (elem < vet[meio])
9             fim = meio - 1;
10        else if (elem > vet[meio])
11            fim = meio + 1;
12        else
13            return meio;
14    }
15    return -1;
16 }
```

- Qual a complexidade desse algoritmo?

Repetições	Tamanho do problema
1	n
2	$n/2$
3	$n/4$
...	...
$\log n$	1

- Assim, são necessárias $\log n$ repetições.
- Como fazemos um número constante de comparações a cada ciclo (duas), podemos concluir que a ordem de complexidade desse algoritmo $O(\log n)$.

- O algoritmo de busca binária tem bom desempenho computacional e deve ser usado sempre que temos dados ordenados armazenados em um vetor.
- Contudo, se precisarmos inserir e remover elementos da estrutura e ao mesmo tempo dar suporte a funções de busca eficientes, a estrutura de vetor não se mostra adequada.
- Precisamos de uma estrutura dinâmica e eficiente: Árvore Binária de Busca - $O(\log n)$.

- Uma Árvore Binária de Busca é tal que ou a árvore é dita vazia ou seu nó raiz contém uma chave e:
 - Todas as chaves da sub-árvore esquerda são menores que a chave da raiz.
 - Todas as chaves da sub-árvore direita são maiores que a chave raiz.
 - As sub-árvore direita e esquerda são também Árvore Binária de Busca.

- Vamos analisar as funções para árvores binária de busca:
 - busca: função que busca um elemento na árvore;
 - insere: função que insere um novo elemento na árvore;
 - retira: função que retira um elemento da árvore.

Operação de Busca

```
1 Arv *abb_busca(Arv *r, int v){  
2     if (r == NULL) return NULL;  
3     else if (r->info > v) return abb_busca(r->esq, v);  
4     else if (r->info < v) return abb_busca(r->dir, v);  
5     else return r;  
6 }
```

Operação de Inserção

```
1 Arv *abb_inserere(Arv *a, int v){
2     if (a == NULL){
3         a = (Arv *) malloc (sizeof(Arv));
4         a->info = v;
5         a->esq = a->dir = NULL;
6     }
7     else if (v < a->info)
8         a->esq = abb_inserere(a->esq, v);
9     else
10        a->dir = abb_inserere(a->dir, v);
11    return a;
12 }
```

- A operação de remoção é mais complexa que a inserção.
- Existem três situações possíveis:
 - Retirar uma raiz que é um nó folha.
 - Retirar uma raiz que tem apenas um filho.
 - Retirar uma raiz que tem dois filhos.

Operação de Remoção

```
1  Arv *abb_retira(Arv *r, int v){
2      if (r == NULL)
3          return NULL;
4      else if (r->info > v)
5          r->esq = abb_retira(r->esq, v);
6      else if (r->info < v)
7          r->dir = abb_retira(r->dir, v);
8      else{
9          //elemento sem filhos
10         if (r->esq == NULL && r->dir == NULL)
11             free(r);
12             r = NULL;
13         //elemento tem filho à direita
14         else if (r->esq == NULL){
15             Arv *t = r;
16             r = r->dir;
17             free(t);
18         }
19         //elemento tem filho à esquerda
20         else if (r->dir == NULL){
21             Arv *t = r;
22             r = r->esq;
23             free(t);
24         }
25         //...
```

Operação de Remoção

```
1 //...
2 //elemento tem dois filhos
3 else{
4     Arv *f = r->esq;
5     while (f->dir != NULL){
6         f = f->dir;
7     }
8     r->info = f->info;
9     f->info = v;
10    r->esq = abb_retira(r->esq, v);
11 }
12 }
13 return r;
14 }
```


- É fácil prever que, após várias operações de inserção/remoção, a árvore tende a ficar desbalanceada.
- Para manter o balanceamento da árvore, devemos tentar manter a altura da árvore no mínimo.
- Algoritmos para manter a árvore balanceada são mais complexos, como já vimos!

Ordenação

- Definição do problema: Dada uma coleção de n elementos, representada em um vetor de 0 a $n-1$, deseja-se obter uma outra coleção, cujos elementos estejam ordenados segundo algum critério de comparação entre os elementos.
- Tipos de algoritmos que serão discutidos: Iterativos e Recursivos
- Alguns exemplos: InsertionSort, **BubbleSort**, RadixSort, HeapSort, **MergeSort**, **QuickSort**.

- Os elementos do vetor devem ser trocados entre si para que fiquem na ordem desejada.

```
1 void bubbleSort(int *p, int t){
2     int i, j, aux;
3     for ( i=0; i < t; i++){
4         for ( j=t-1; j > i; j-- ){
5             if ( p[j] < p[j-1]){
6                 aux = p[j];
7                 p[j] = p[j-1];
8                 p[j-1] = aux;
9             }
10        }
11    }
12 }
```

<http://www.youtube.com/watch?v=lyZQPjUT5B4>

- A idéia básica do MergeSort é criar uma sequência ordenada a partir de duas outras também ordenadas. Para isso, o algoritmo MergeSort divide a sequência original em pares de dados, agrupa estes pares na ordem desejada; depois agrupa as sequências de pares já ordenados, formando uma nova sequência ordenada de quatro elementos, e assim por diante, até ter toda a sequência ordenada.

- Os três passos úteis dos algoritmos dividir-para-conquistar, que se aplicam ao MergeSort são:
 - **Dividir:** Dividir os dados em subsequências pequenas; Este passo é realizado recursivamente, iniciando com a divisão do vetor de n elementos em duas metades, cada uma das metades é novamente dividida em duas novas metades e assim por diante, até que não seja mais possível a divisão (ou seja, sobrem n vetores com um elemento cada).
 - **Conquistar:** Classificar as duas metades recursivamente aplicando o mergesort;
 - **Combinar:** Juntar as duas metades em um único conjunto já classificado. Para completar a ordenação do vetor original de n elementos, faz-se o merge ou a fusão dos sub-vetores já ordenados.

http://www.youtube.com/watch?v=XaqR3G_NVoo

MergeSort

```
1 void merge(int a[], int low, int
    high, int mid){
2     int i, j, k, c[50];
3     i=low;
4     j=mid+1;
5     k=low;
6     while((i<=mid)&&(j<=high)){
7         if(a[i]<a[j]){
8             c[k]=a[i];
9             k++;
10            i++;
11        }
12        else{
13            c[k]=a[j];
14            k++;
15            j++;
16        }
17    }
18    while(i<=mid){
19        c[k]=a[i];
20        k++;
21        i++;
22    }
23    while(j<=high){
24        c[k]=a[j];
25        k++;
26        j++;
27    }
28    for(i=low; i<k; i++)
29        a[i]=c[i];
30 }
```

```
1 void mergeSort(int a[], int low, int
    high){
2     int mid;
3     if(low>high){
4         mid=(low+high)/2;
5         mergeSort(a, low, mid);
6         mergeSort(a, mid+1, high);
7         merge(a, low, high, mid);
8     }
9 }
```

- Determina-se um elemento pivô. O pivô é posicionado dentro do vetor de tal forma que, todos à esquerda do pivô são menores que ele e, todos à direita do pivô são maiores. O pivô "divide" o vetor em dois subvetores. Recursivamente o quicksort é realizado na primeira metade do vetor e na segunda metade.

- Algoritmo: Seja x o vetor a ser ordenado e n o número de elementos de x . Seja a um elemento de x escolhido ao acaso (por exemplo, $a=x[0]$). Suponha que os elementos de x estejam divididos de tal forma que a é colocado na posição j e as seguintes condições são verdadeiras:
 - Todos os elementos nas posições de 0 a $j-1$ são menores que a .
 - Todos os elementos nas posições de $j+1$ a $n-1$ são maiores ou iguais a a .
- Então a está na posição correta no vetor. Se este processo for repetido para os sub-vetores $x[0]$ a $x[j-1]$ e $x[j+1]$ a $x[n-1]$, o resultado é o vetor ordenado.

<http://www.youtube.com/watch?v=kDgvnbUIqT4>

QuickSort

```
1  int partition( int a[], int l, int r) {
2      int pivot, i, j, aux;
3      pivot = a[l];
4      i = l;
5      j = r+1;
6      while(1){
7          do ++i; while( a[i] <= pivot && i <= r );
8          do --j; while( a[j] > pivot );
9          if( i >= j ) break;
10         aux = a[i];
11         a[i] = a[j];
12         a[j] = aux;
13     }
14     aux = a[l];
15     a[l] = a[j];
16     a[j] = aux;
17     return j;
18 }
19
20 void quickSort( int a[], int l, int r){
21     int j;
22     if( l > r ){
23         //dividir e conquistar
24         j = partition( a, l, r);
25         quickSort( a, l, j-1);
26         quickSort( a, j+1, r);
27     }
28 }
```

No main o QuickSort será chamado como quickSort(vetor, 0, tam);

Hash/Dispersão

- Considerar o problema de pesquisar um determinado valor em um vetor.
 - Se o vetor não está ordenado, a pesquisa requer $O(n)$ de complexidade.
 - Se o vetor está ordenado, pode-se fazer a pesquisa binária que requer uma complexidade de $O(\log n)$.
 - Não parece haver melhor maneira de resolver o problema com custos melhorados.

- Ideia: Tabelas Hash/Dispersão:
 - Poderia haver uma maneira de resolver o problema em $O(1)$.
 - Se o vetor estiver ordenado de uma determinada maneira.
- Solução?
 - Arranjar uma função mágica que, dado um determinado valor a pesquisar, nos diga exatamente a posição no vetor.
 - Esta é uma função chamada **função de dispersão/hash**.

- As **tabelas de hash** são um tipo de estruturação criado para o armazenamento de informação, e são uma forma extremamente simples, fácil de se implementar e, intuitiva para se organizar grandes quantidades de dados.
- Possui como ideia central a divisão do universo de dados a ser organizado em subconjuntos mais facilmente gerenciáveis.
- **Permitir armazenar e procurar rapidamente grande quantidade de dados**

- As tabelas de hash são constituídas por 2 conceitos fundamentais:
 - *Tabela de Hash*: estrutura que permite o acesso aos subconjuntos.
 - *Função de Hash*: função que realiza um mapeamento entre os valores de chaves e as entradas na tabela.

- Criar um critério simples para dividir este universo em subconjuntos com base em alguma qualidade do domínio das chaves.
 - Possuir um índice que permita encontrar o início do subconjunto certo, depois de calcular o valor de hash.
 - Isto é uma tabela de hash.

- Saber em qual subconjunto procurar e colocar uma chave.
 - indicar quantos subconjuntos se pretende.
 - criar uma regra de cálculo que, dada uma chave, determine em que subconjunto se deve procurar pelos dados com esta chave ou colocar estes dados (caso seja um novo elemento).
 - Isto é chamado função de hash.

- Gerir estes subconjuntos bem menores com um métodos simples:
 - Possuir uma estrutura ou um conjunto de estruturas de dados para os subconjuntos.
 - Existem duas filosofias: hashing fechado (ou de endereçamento aberto) ou o hashing aberto (ou encadeado).
- Alguns dos problemas que se colocam quando usamos tabelas de hash são:
 - determinar uma função de hash que minimize o número de colisões;
 - obter os mecanismos eficientes para tratar as colisões.

Funções de Hash - Exemplos

- Números x (chave) de 0 a 99 (dois dígitos)
- tam tamanho da tabela
- Pode-se construir uma função que coloque x no vetor em termos do algarismo das dezenas.
 $f = x / tam$ ($/$ = divisão inteira)
- Ou pode-se construir construir uma função que coloque x no vetor em termos do algarismo das unidades.
 $f = x \% tam$ ($\%$ = resto da divisão)

Funções de Hash - Exemplos

```
1 int hash (int key, int tam) {  
2     return (key % tam);  
3 }
```

- Objetivos das funções de Hash:
 - deve ser eficiente.
 - deve distribuir todos os elementos uniformemente por todas as posições da tabela.

Método normal

```
1 int hash (char key[], int tam) {  
2     int valHash = 0, k = 0;  
3     while (key[k] != '\0') {  
4         valHash = valHash + int(key[k]);  
5         k++;  
6     }  
7     return (valHash % tam);  
8 }
```

Funções de Hash - Exemplos Strings

Usando a regra de Horner:

```
1 int hash (char key[], int tam) {  
2     int valHash = 0, a = 127, k = 0;  
3     while (key[k] != '\0') {  
4         valHash = valHash * a + int(key[k]);  
5         k++;  
6     }  
7     return (valHash % tam);  
8 }
```

- Quando a função de Hash é imperfeita poderão existir dois valores a serem colocados na mesma posição do vetor. A isto chama-se uma colisão.
- As colisões são normalmente tratadas como quem chega primeiro serve-se; isto é, o primeiro elemento a chegar a uma posição fica com ela.
- Terá de se arranjar uma solução eficiente para se determinar o que se deve fazer com o segundo valor que deveria ser colocado na mesma posição.

- O que fazer quando dois valores diferentes tentam ocupar a mesma posição?
 - Solução 1 (**hashing fechado** ou **endereçamento aberto**): procura a partir dessa posição uma posição vazia.
 - Solução 2: usar uma segunda função de hash, uma terceira, uma quarta, ?
 - Solução 3 (**hashing aberto** ou **encadeamento separado**): usa a posição do vetor como cabeça (head) de uma lista que vai conter todas as colisões dessa posição.

Hashing Fechado

- Cria-se uma estrutura que assinale cada posição:
 - Livre (vazia e nunca ocupada)
 - Ocupada
 - Removida (vazia mas já ocupada)
- Pesquisar/remover:
 - Termina quando encontra o objeto
 - Termina quando encontra uma posição livre
 - Continua a pesquisa quando encontra uma posição removida
- Inserir (insere o objeto quando encontra):
 - uma posição assinalada como livre ou removida.

Hashing Aberto

- A estruturação dos dados é talvez a forma mais intuitiva de se implementar o conceito de Hashing.
- Consiste em ter um vetor de apontadores, com dimensão N , em que cada elemento do vetor contém uma ligação para uma lista dos elementos a guardar.
- A pesquisa de um elemento efetua-se da seguinte forma:
 - A partir de uma chave, calcular qual o elemento do vetor é a cabeça da lista que se pretende.
 - Usando um qualquer algoritmo para o efeito, pesquisar o elemento dentro daquela lista.

Tabela Hash/Dispersão - Hashing Aberto

Exemplo 1:

- Armazenar os elementos: 19, 26, 33, 70, 79, 103 e 110.
- Usando a função de hash: $\text{hash}(x) = x \% 10$

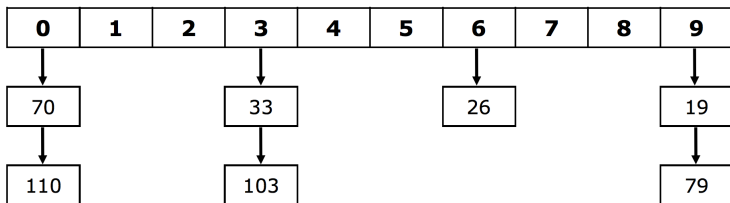


Tabela Hash/Dispersão - Hashing Aberto

Exemplo 2:

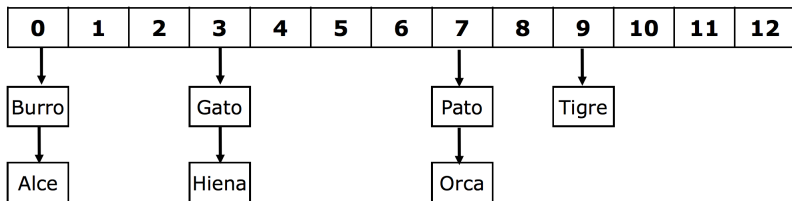
- Cada elemento do vetor é um apontador para uma lista de estruturas com o mesmo valor da função de hash.
- Supondo que o vetor tem tamanho $\text{tam} = 13$ e que a função de hash é a que consta na tabela que se segue

Chave	A,B	C,D	E,F	G,H	I,J	K,L	M,N	O,P	Q,R	S,T	U,V	X,Y	W,Z
Hash	0	1	2	3	4	5	6	7	8	9	10	11	12

Tabela Hash/Dispersão - Hashing Aberto

Exemplo 2:

- Armazenar os elementos: Alce, Burro, Gato, Hiena, Pato, Orca e Tigre.



- Reduz o número de comparações, quando comparado com a pesquisa sequencial.
- Necessidade de espaço em memória para armazenar o vetor de listas.

- Dispersão aberta (DA) vs. Fechada (DF)
 - Suporta a primitiva de remoção.
 - As listas de colisão não se cruzam.
 - O erro por defeito do pré-dimensionamento não é tão grave.
 - Se a tabela estiver em arquivo poupam-se muitos acessos com a sondagem linear, porque em geral os elementos que colidem estão fisicamente próximos.

Tabela Hash/Dispersão - Vantagens e Desvantagens

- Vantagens da Dispersão:
 - A complexidade das primitivas suportadas (pesquisa, inserção e remoção na DA), no caso esperado é constante.
 - A técnica é eficiente e é só depende do fator de ocupação e da qualidade das funções (na DF).
- Problemas da Dispersão:
 - Não é uma estrutura dinâmica (o redimensionamento é necessário)
 - Não suporta primitivas que se baseiam em relações de ordem dos elementos (mínimo, máximo, percurso ordenado).
 - A complexidade das primitivas suportadas (pesquisa, inserção e remoção na DA), no pior caso é linear no nº de elementos da tabela.

- Comentamos ao longo da aula sobre complexidade, algoritmo eficiente, $O(n)$... Mas o que significa isso?

- 1. Crie uma Agenda de Contatos (Nome, Telefone, Endereço) com uma Busca Binária baseada no nome do contato.
- 2. Crie uma Agenda de Contatos (Nome, Telefone, Endereço) baseada em Tabela Hash.

Na próxima aula...

Exercícios para prova