

Trabalho Pratico II - Compressor de Arquivos

Joao Paulo Mendes de Sa

1 Introdução

O algoritmo de Huffman é capaz de gerar um código de tamanho variável para um conjunto de símbolos com a finalidade de otimizar a quantidade de bits necessárias para armazenar o conteúdo de um arquivo. Códigos menores são usados para representar símbolos que ocorrem no arquivo com mais frequência e os maiores para símbolos com menor frequência. O algoritmo de Huffman gera uma árvore binária onde o caminho até cada folha é o seu código. Dessa maneira é impossível gerar um código que seja o prefixo de outro código na mesma árvore pois apenas folhas guardam símbolos. Os passos básicos do algoritmo de Huffman para gerar a árvore sabendo a frequência de cada símbolo são:

1. Crie uma fila de prioridades com nodos de uma árvore onde os itens de menor frequência tendo maior prioridade.
2. Enquanto há mais de um nodo na fila:
 - (a) Remova os dois primeiros nodos da fila
 - (b) Crie um novo nodo e faça a sua esquerda e direita apontar para os dois nodos retirados
 - (c) A frequência desse nodo é igual a soma das frequências dos nodos retirados
 - (d) Insira esse nodo de novo na fila
3. O nodo que resta é a raiz da árvore

Para determinar o código para um símbolo basta navegar da raiz até o nodo que contém o símbolo. Cada esquerda é um 0 e cada direita é um 1. A distância de uma folha até a raiz é o tamanho do código gerado.

2 Implementação

2.1 Estrutura de Dados

Foi criado um tipo abstrato de dados para manipular o heap, um para gerar a árvore de Huffman, e outro para comprimir e descomprimir o arquivo dado uma árvore contendo a codificação.

O heap armazenava todas as folhas na árvore e continha informações sobre o símbolo original, a frequência que ele aparecia, o tamanho do código gerado pelo algoritmo de Huffman, o código gerado, e ponteiros para outras folhas. A árvore era composta desses mesmos subitens do heap mais um ponteiro adicional para a raiz da árvore.

Para implementar a fila de prioridade foram usados dois heaps ao invés de uma lista encadeada. Dessa maneira o algoritmo é muito mais eficiente. Com apenas uma lista ordenada remoção seria $O(1)$ pois bastaria fazer a cabeça da lista apontar para o próximo, inserção seria $O(n)$ pois seria necessário buscar pela lista inteira até achar o lugar de inserção, e ordenação $O(n^2)$ através de seleção ou inserção. Usando dois heaps remoção virou $O(\log(n))$ logo que a

função `heapify()` mantém as propriedades de uma heap navegando ele como uma árvore binária, inserção $O(1)$ pois como a soma de dois nodos sempre seria maior que todos nodos anteriores bastava reinserir o nodo na ultima posição do segundo heap, e ordenação $O(n\log(n))$ usando heapsort.

2.2 Funções e Procedimentos

```
unsigned long file_len(FILE* file);
```

Recebe um arquivo e retorna o seu tamanho em bits.

```
void heap_init(heap_t *heap);
```

Aloca memoria para armazenar os símbolos no heap e inicializa todos eles.

```
void heap_dealloc(heap_t *heap);
```

Libera a memoria alocada ao heap.

```
void heap_fill(heap_t *heap, FILE *file);
```

Ler um arquivo e determina a frequência de cada simbolo no arquivo

```
void heap_condense(heap_t *heap);
```

Apos o `heap_fill()` essa função remove do vetor do heap os items que tem frequência igual a 0 e determina o numero de símbolos diferentes que apareceram

```
int heapCmp_freq(byte_t a, byte_t b);  
int heapCmp_symb(byte_t a, byte_t b);
```

Ambas usados com ponteiros para funções pelas funções que mantem o heap. A primeira ordena por frequência do simbolo e a segunda por ordem alphabetic do simbolo.

```
void heapify(heap_t *heap, int father, int (*cmp)(byte_t, byte_t));
```

Função navega o vetor do heap como uma árvore e verifica se o heap ainda é um heap. Se não o concerta.

```
void heap_build(heap_t *heap, int (*cmp)(byte_t, byte_t));
```

Cria um heap a partir de um vetor ordenado. O critério e decido por uma outra função passada por ponteiro.

```
void heap_sort(heap_t *heap, int (*cmp)(byte_t, byte_t));
```

Ordena um vetor usando heapsort e o critério dado pela função. Primeiro cria um heap do vetor com `heap_build()` e depois remove primeiro item e o coloca na ultima posição do vetor. Conserta o heap com `heapify()` e repete ate o vetor estar ordenado.

```
void heap_clone(heap_t original, heap_t* clone);
```

Duplica um heap.

```
byte_t heap_extract_min(heap_t *heap, int(*cmp)(byte_t, byte_t));
```

Remove of menor item to heap e chama heapify() para consertar depois.

```
byte_t* huff_build_tree(heap_t huffQueue);
```

Implementa o algoritmo de Huffman. Duplica o heap e cria um heap auxiliar. Enquanto houver mais de um elemento entre os dois heap remove os dois com menores frequências e cria uma subárvore e o reinsere no fim do segundo heap. Retorna a raiz para a árvore gerada.

```
void tree_dealloc(byte_t *node);
```

Faz um caminhamento preordem da árvore e desaloca a memória usada por cada nodo.

```
int bsearchCmp_symb(const void* a, const void* b);
```

Usada por bsearch() como critério para buscar um símbolo.

```
void bin_to_dec(char* binary, unsigned short *dec, short bits);
```

Pega um número binário representado por uma string (little endian) e converte para base 10

```
void huff_build_decode_table(heap_t *heap, tree_t tree);
```

Navega árvore de Huffman e determina o código usado para representar cada símbolo. Escreve essa informação no heap.

```
void file_write_header(FILE *outputFile, heap_t heap,  
                      unsigned long inputFileLen);
```

Escreve of cabeçalho no início do arquivo a ser comprimido.

```
%void dec_to_bin(unsigned short dec, char* prefixBin, short bits);
```

Faz of inverso de bin_to_dec().

```
void write_bit(FILE *outputFile, short prefix, short prefixLen);
```

Usando bitmasking escreve cada bit do código de Huffman de um símbolo a um buffer. Quando o buffer é preenchido por oito bits escreve o byte no arquivo.

```
void file_compress(FILE *outputFile, FILE *inputFile, heap_t heap,  
                  unsigned long inputFileLen);
```

Efetua a compressão. Ler o arquivo a ser comprimido e produz um bitstream usando write_bit() para escrever no arquivo a ser comprimido.

```
unsigned long file_parse_header(FILE *compressedFile, tree_t *huff);
```

Ler o cabeçalho de um arquivo e reconstrói a árvore de Huffman. Começa por um nodo raiz e vai navegando pelo código de Huffman de cada character. Se um nodo não existe então a função aloca espaço para ele até chegar na folha correta. Repete isso para todos símbolos.

```
void file_decompress(FILE *compressedFile, FILE *outputFile, tree_t huff,  
                    unsigned long fileLen);
```

Efetua a descompressão do bitstream. Usando a árvore reconstituída ele lê um byte e escreve no buffer. Esse buffer é convertido para um número binário e o caminhamento da árvore é feito até encontrar uma folha. O símbolo decodificado é escrito no arquivo. Isso repete até o buffer esvaziar. Quando vazio ele lê outro bit até o arquivo descomprimido esteja do seu tamanho original.

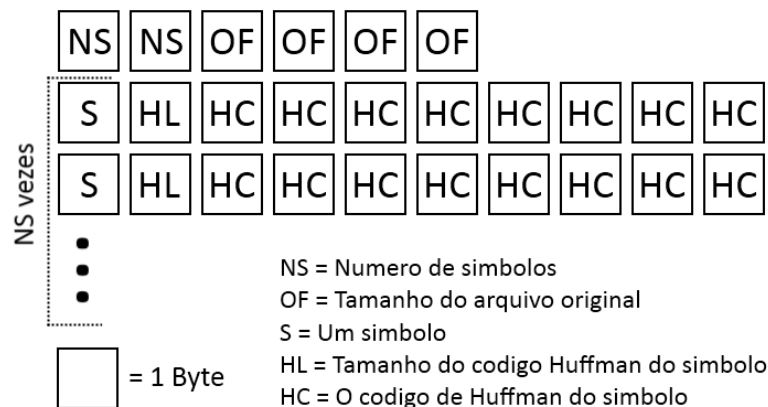
2.3 Programa Principal

é necessário ler o arquivo a ser comprimido duas vezes. Uma para determinar os símbolos que aparecem que nessa caso foi cada byte e a frequência na qual ele ocorre. A segunda leitura é a que realmente comprime o arquivo usando o código gerado pelo algoritmo de Huffman. Primeiro é escrito um cabeçalho no arquivo comprimido com as informações necessárias para decodificá-lo e finalmente o conteúdo do arquivo a ser comprimido se transforma no bitstream codificado que é escrito no arquivo comprimido. Para descomprimir o arquivo bastava apenas reconstruir a árvore. Um processo simples sendo que todas informações necessárias estão presente no cabeçalho.

2.4 Organização do Código, Decisões de Implementação e Detalhes Técnicos

O código está dividido em quatro arquivos. `huff.c` e `huff.h` contêm todas as funções necessárias para efetuar a compressão e descompressão de um arquivo. `comprima.c` contém o procedimento para comprimir um arquivo e `descomprima.c` faz o inverso. Para significar a implementação do algoritmo o maior código de Huffman que o programa é capaz de gerar é de 64 bits pela limitação que o maior tipo em ANSI C90 é `unsigned long` logo isso é menor que o limite teórico de 256 bits. Para retificar essa limitação eu podia separar o código de um símbolo em 4 variáveis que daria 256 bits no total. Mas essa alteração aumentaria a complexidade do código significativamente. Para verificar se esse limite era encontrado em casos reais eu gerei 10 arquivos com conteúdos randômicos entre 100MB e 1GB e nunca encontrei um caso onde essa limitação ficou aparente.

O formato usado no cabeçalho de decodificação foi:



O compilador usado foi `gcc 4.6.020110429` no sistema operacional Arch Linux 2.6.38-ARCH. Para executar o compressor basta digitar na linha de commando `./comprima FILEOUT FILEIN` e para o descompressor `./descomprima FILEOUT FILEIN`.

3 Analise de Complexidade

A variável x é definido como o tamanho do arquivo em bits. A variável n é definido como o numero de símbolos diferentes (no máximo 256).

```
unsigned long file_len(FILE* file);
```

Não depende da entrada logo determina o tamanho em $O(1)$.

```
void heap_init(heap_t *heap);
```

Não depende da entrada longo aloca e inicializa em $O(1)$.

```
void heap_dealloc(heap_t *heap);
```

Não depende da entrada long desaloca em $O(1)$.

```
void heap_fill(heap_t *heap, FILE *file);
```

Para cada byte do arquivo incrementa a frequência de um simbolo no heap logo $O(x)$.

```
void heap_condense(heap_t *heap);
```

Os dois loops aninhados tem tamanho fixo logo não depende da entrada então $O(1)$.

```
int heapCmp_freq(byte_t a, byte_t b);  
int heapCmp_symb(byte_t a, byte_t b);
```

Cada um faz apenas uma comparação independente do tamanho logo $O(1)$.

```
void heapify(heap_t *heap, int father, int (*cmp)(byte_t, byte_t));
```

Navega um vetor como uma arvore, no pior caso navega de uma raiz ate uma folha que é $O(\log(n))$.

```
void heap_build(heap_t *heap, int (*cmp)(byte_t, byte_t));
```

Chama heapify() para cada metade de n longo $O(n\log(n))$.

```
void heap_sort(heap_t *heap, int (*cmp)(byte_t, byte_t));
```

Constrói o heap com heap_build() e depois para cada n a função remove a raiz do heap e reinsere na ultima posição do vetor. $\max(O(n\log(n)), O(n\log(n))) = O(n\log(n))$.

```
void heap_clone(heap_t original, heap_t* clone);
```

Um loop depende no numero de símbolos logo $O(n)$.

```
byte_t heap_extract_min(heap_t *heap, int(*cmp)(byte_t, byte_t));
```

Chama heapify() uma vez então $O(\log(n))$.

```
byte_t* huff_build_tree(heap_t huffQueue);
```

Um loop para inicializar um heap auxiliar $O(n)$. Um loop para o algoritmo de Huffman que também executa n vezes e chama heap_extract_min() duas vezes logo é $2O(n)O(\log(n)) = O(n\log(n))$.

```
void tree_dealloc(byte_t *node);
```

Navega todos os nodos da arvore que são no máximo $2n - 1$ logo $O(n)$.

```
int bsearchCmp_symb(const void* a, const void* b);
```

Uma comparação $O(1)$.

```
void bin_to_dec(char* binary, unsigned short *dec, short bits);
```

Depende do tamanho de bits que é sempre 8 ou 16 logo $O(1)$.

```
void huff_build_decode_table(heap_t *heap, tree_t tree);
```

Navega todos os nodos da arvore $(2n - 1)$ em preordem $O(n)$. Nas folhas faz uma busca binaria $O(\log(n))$, Escreve o código em $O(1)$. $O(n)O(\log(n))O(1) = O(n\log(n))$.

```
void file_write_header(FILE *outputFile, heap_t heap,
    unsigned long inputFileLen);
```

O tamanho do cabeçalho é determinado pelo numero de símbolos $O(n)$.

```
%void dec_to_bin(unsigned short dec, char* prefixBin, short bits);
```

Depende do tamanho de bits que é sempre 8 ou 16 logo $O(1)$.

```
void write_bit(FILE *outputFile, short prefix, short prefixLen);
```

Depende do tamanho do código de um simbolo que no pior caso de uma arvore completamente desbalanceada de altura n seria então $O(n)$.

```
void file_compress(FILE *outputFile, FILE *inputFile, heap_t heap,
    unsigned long inputFileLen);
```

Para cada bit do tamanho do arquivo faz um busca binaria $O(\log(n))$ e chama `write_bit()` que é $O(n)$. $O(x)\max(O(n), O(\log(n))) = O(x)O(n)$.

```
unsigned long file_parse_header(FILE *compressedFile, tree_t *huff);
```

O tamanho do cabeçalho é n e para cada simbolo ela navega arvore dependendo do tamanho do código do simbolo que no pior caso igual a n então $O(n^2)$

```
void file_decompress(FILE *compressedFile, FILE *outputFile, tree_t huff,
    unsigned long fileLen);
```

Depende do tamanho original do arquivo long $O(x)$.

3.1 Programa principal

A complexidade do programa de compressão então é

$$O(n)O(x) = \max(O(1), O(x), O(1), O(n\log(n)), \\ O(n\log(n)), O(n\log(n)), O(n\log(n)), \\ O(x)O(n), O(1))$$

A complexidade do programa de descompressão é

$$O(n^2) + O(x) = \max(O(n^2), O(x))$$

Como no pior caso $n = 256$ podemos considerar $O(256)$ e $O(256^2)$ sendo igual $O(1)$ o que deixa a complexidade de compressão e descompressão sendo $O(x)$.

4 Testes

Foram feitos 10 testes em arquivos diferentes and conferido o checksum do arquivo antes de compressão e depois de compressão e descompressão

Tipo	T Original (bytes)	T Comprimido (bytes)	Compressão
pdf	4315810	4079264	94.5 %
executável	23461	19325	82.4 %
vazio	0	0	N/A
jpg	944787	946291	100.2 %
png	2027978	2030544	100.1 %
png	11120	11791	106.0 %
txt	1650	938	58.7 %
txt	109925833	70584126	64.2 %
mkv	203499149	203501715	100.0 %
avi	547350528	545511673	99.7 %

```
❶ ❷ ❸ ❹ ❺ []= urxvt VOL: 69% | CHR: 95% | 2011/05/31 01:01
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/book.com tests/book.pdf; ./descomprima tests/book0.pdf tests/book.com; md5sum tests/book.pdf tests/book0.pdf
af502ba10787844979f9583262b2aae7 tests/book.pdf
af502ba10787844979f9583262b2aae7 tests/book0.pdf
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/comprima.com tests/comprima; ./descomprima tests/comprima0 tests/comprima.com; md5sum tests/comprima tests/comprima0
c16000eb08020b7be89cccb7a75d3d37 tests/comprima
c16000eb08020b7be89cccb7a75d3d37 tests/comprima0
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/empty.com tests/empty; ./descomprima tests/empty0 tests/empty.com; md5sum tests/empty tests/empty0
d41d8cd98f00b204e9800998ecf8427e tests/empty
d41d8cd98f00b204e9800998ecf8427e tests/empty0
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/image1.com tests/image1.jpg; ./descomprima tests/image10.jpg tests/image1.com; md5sum tests/image1.jpg tests/image10.jpg
c5541e4fc58fa7bac673a69f7a5a0058 tests/image1.jpg
c5541e4fc58fa7bac673a69f7a5a0058 tests/image10.jpg
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/image2.com tests/image2.png; ./descomprima tests/image20.png tests/image2.com; md5sum tests/image2.png tests/image20.png
662f31e3dfdfa14931b6a8471ae8b657 tests/image2.png
662f31e3dfdfa14931b6a8471ae8b657 tests/image20.png
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/image3.com tests/image3.png; ./descomprima tests/image30.png tests/image3.com; md5sum tests/image3.png tests/image30.png
e9e8badaaaf06f8b0e879c1798c554d6 tests/image3.png
e9e8badaaaf06f8b0e879c1798c554d6 tests/image30.png
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/text1.com tests/text1.txt; ./descomprima tests/text10.txt tests/text1.com; md5sum tests/text1.txt tests/text10.txt
2fa8f4ec4d59d87c44d4d5dfcebf104e6 tests/text1.txt
2fa8f4ec4d59d87c44d4d5dfcebf104e6 tests/text10.txt
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/text2.com tests/text2.txt; ./descomprima tests/text20.txt tests/text2.com; md5sum tests/text2.txt tests/text20.txt
8e56012b762df613a82698f2ad5a37f2 tests/text2.txt
8e56012b762df613a82698f2ad5a37f2 tests/text20.txt
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/video1.com tests/video1.mkv; ./descomprima tests/video10.mkv tests/video1.com; md5sum tests/video1.mkv tests/video10.mkv
782621072f06404e69b7a54e6b1f63f1 tests/video1.mkv
782621072f06404e69b7a54e6b1f63f1 tests/video10.mkv
[kakiray@SUSSMHN filearchiver]$ ./comprima tests/video2.com tests/video2.avi; ./descomprima tests/video20.avi tests/video2.com; md5sum tests/video2.avi tests/video20.avi
c43915b97dd375f560f5b85e8f146740 tests/video2.avi
c43915b97dd375f560f5b85e8f146740 tests/video20.avi
[kakiray@SUSSMHN filearchiver]$ █

❶ ❷ ❸ ❹ ❺ []= urxvt VOL: 69% | CHR: 96% | 2011/05/31 01:03
[kakiray@SUSSMHN filearchiver]$ ls -l tests/
total 2502884
-rw-r--r-- 1 kakiray users 4079264 May 31 00:12 book.com
-rw-r--r-- 1 kakiray users 4315810 May 31 00:12 book0.pdf
-rw-r--r-- 1 kakiray users 4315810 Oct 18 2010 book.pdf
-rwxr-xr-x 1 kakiray users 23461 May 31 00:03 comprima
-rw-r--r-- 1 kakiray users 19325 May 31 00:14 comprima.com
-rw-r--r-- 1 kakiray users 23461 May 31 00:14 comprima0
-rw-r--r-- 1 kakiray users 0 May 31 00:05 empty
-rw-r--r-- 1 kakiray users 0 May 31 00:21 empty.com
-rw-r--r-- 1 kakiray users 0 May 31 00:21 empty0
-rw-r--r-- 1 kakiray users 944787 May 31 00:23 image1.com
-rw-r--r-- 1 kakiray users 944787 May 31 00:23 image10.jpg
-rwxr-xr-x 1 kakiray users 944787 Jul 10 2008 image1.jpg
-rw-r--r-- 1 kakiray users 2030544 May 31 00:24 image2.com
-rw-r--r-- 1 kakiray users 2027978 May 31 00:24 image20.png
-rwxr-xr-x 1 kakiray users 2027978 Jan 13 23:36 image2.png
-rw-r--r-- 1 kakiray users 11791 May 31 00:25 image3.com
-rw-r--r-- 1 kakiray users 11120 May 31 00:25 image30.png
-rw-r--r-- 1 kakiray users 11120 May 9 17:01 image3.png
-rw-r--r-- 1 kakiray users 938 May 31 00:26 text1.com
-rw-r--r-- 1 kakiray users 1650 May 31 00:26 text10.txt
-rw-r--r-- 1 kakiray users 1650 May 29 12:44 text1.txt
-rw-r--r-- 1 kakiray users 70584126 May 31 00:28 text2.com
-rw-r--r-- 1 kakiray users 109925833 May 31 00:29 text20.txt
-rw-r--r-- 1 kakiray users 109925833 May 3 22:53 text2.txt
-rw-r--r-- 1 kakiray users 203501715 May 31 00:33 video1.com
-rw-r--r-- 1 kakiray users 203499149 May 31 00:38 video10.mkv
-rwxr-xr-x 1 kakiray users 203499149 Jan 13 20:26 video1.mkv
-rwxr-xr-x 1 kakiray users 547358528 Oct 27 2006 video2.avi
-rw-r--r-- 1 kakiray users 545511673 May 31 00:45 video2.com
-rw-r--r-- 1 kakiray users 547358528 May 31 00:59 video20.avi
[kakiray@SUSSMHN filearchiver]$ █
```

5 Conclusão

Em arquivos muito grandes o numero de símbolos são quase sempre chega a 256 que torna n irrelevante logo programa é limitado apenas pelo tamanho da entrada. Em alguns arquivos a compressão resultou em um arquivo de tamanho maior que o original. Esses arquivos já tiveram alguma forma de compressão aplicadas a eles. Para verificar essa hipótese eu testei comprimir um arquivo anteriormente comprimido e ele também acabou tendo um tamanho maior que o original. A implementação ocorreu sem muitos problemas alem do ocasional segmentation fault.

6 Referencias

- Cormen T., Leiserson C., Rivest R., Stein C., Introduction to Algorithms 3rd Edition 2009.
- Case S. <http://www.programmersheaven.com/download/2253/0/ZipView.aspx>. 1991.
- http://en.wikipedia.org/Huffman_coding