

Politechnika Śląska
Wydział Automatyki, Elektroniki i Informatyki

Podstawy Programowania Komputerów

Temat:
Kompresja plików metodą kodowania Huffmana

Autor	Piotr Hasiec
Prowadzący	Dr Inż. Bożena Wieczorek
Rok akademicki	2019/2020
Kierunek	Informatyka
Rodzaj studiów	SSI
Semestr	2
Sekcja	12
Termin oddania sprawozdania	2020-09-09

1. Treść zadania

Napisać program, dokonujący bezstratnej kompresji plików metodą kodowania kodem prefikso-
wym Huffmana.

Program uruchamiany jest z linii poleceń z wykorzystaniem następujących przełączników (kolej-
ność przełączników jest dowolna):

- k jeśli chcemy skompresować plik,
- d jeśli chcemy rozpakować plik,
- i <NAZWA.ROZSZERZENIE> plik wejściowy
- o <NAZWA.ROZSZERZENIE> plik wyjściowy

2. Analiza zadania

2.1. Struktury danych

W programie wykorzystane są struktury drzewiaste. Pierwszą ze struktur jest drzewo będące
wynikiem działania funkcji `node* add_to_tree(node* leafs)`; będące drzewem kodów
Hufmanna, drugą jest struktura kopca typu min zaimplementowana w celu użycia jej jako kolejki
priorytetowej.

2.2. Algorytmy

Użyte w programie algorytmy to algorytm budowy drzewa kodów Huffmana będący algorytmem
zachłannym oraz algorytmy słownikowe kopca typu min.

3. Specyfikacja zewnętrzna

Program jest uruchamiany z linii poleceń. Należy przekazać do programu nazwy plików: wejścio-
wego i wyjściowego po odpowiednich przełącznikach (odpowiednio: -i, -o) oraz procedurę wykony-
waną na tych plikach – kompresję lub rozpakowywanie (odpowiednio: -k, -d)

```
Program.exe -i dane_wejsciove.txt -o wynik -k  
Program.exe -o wyjscie.txt -i dane - d
```

Pliki wejściowy jak i plik wyjściowy może posiadać dowolne rozszerzenie. Należy jednak zwrócić
uwagę żeby przy dekompresji pliku podać odpowiednie rozszerzenie.

Podanie niewłaściwych nazw pliku powoduje otrzymanie komunikatów:

Nie udało się uzyskać dostępu do co najmniej jednego pliku.

W każdym przypadku podania niewłaściwych danych, program kończy swoje działanie.

4. Specyfikacja wewnętrzna

Program został zrealizowany zgodnie z paradygmatem strukturalnym.

4.1. Ogólna struktura programu

W przypadku funkcjonalności kompresji pliku program w pierwszej kolejności zlicza liczbę wystąpień każdego z 256 możliwych bajtów. Złożoność tej procedury wynosi $O(n)$. Następnie na podstawie tych danych budowana jest kolejka priorytetowa typu min, drzewo kodów Huffmana oraz słownik kodów. Złożoność czasowa tych trzech algorytmów wynosi $O(1)$. Nie jest to prawda w ogólności ponieważ: złożoność procedury budowy kopca wynosi: $O(s)$ gdzie s jest liczbą elementów tablicy. Złożoność budowy drzewa Huffmana wynosi $O(s \lg s)$, a procedury `void codematrix(node* A, int** code, int i)` polegającej na przejściu całego drzewa - $O(s)$. Jednak we wszystkich przypadkach $s = 256$ zatem złożoności redukują się do czasu stałego. Ostatnią procedurą funkcjonalności kodowania jest przetłumaczenie nie zakodowanego pliku. Operacja ta wykonuje się w czasie liniowym.

W przypadku funkcjonalności dekompresji pliku program wczytuje tablicę częstości występowania bajtów zapisaną w pliku - $O(1)$, buduje drzewo kodów Huffmana - $O(1)$, wczytuje X bajtów, tłumaczy je na ciąg 0 i 1 zapisanych w przeznaczony do tego tablicy i na jego podstawie dekoduje plik. Dekodowanie pliku odbywa się poprzez n -krotne przejście drzewa Huffmana którego średnia (ważona częstością występowania bitów) głębokość zależy od struktury pierwotnego pliku. Dolnym ograniczeniem wartości n jest liczba bajtów w pierwotnym (nie zakodowanym pliku). W rzeczywistości jest ona większa ponieważ czasem wczytany ciąg 0 i 1 będzie w sobie zawierał niepełne słowo kodowe. Takie zdarzenie jednak może nastąpić jedynie A/X razy (A – liczba bajtów w zakodowanym pliku, X – liczba wczytywanych za jednym razem bajtów). Zatem złożoność procedury dekodowania pliku jest liniowa.

4.2. Szczegółowy opis typów i funkcji

Szczegółowy opis typów i funkcji znajduje się w załączniku na końcu dokumentu.

5. Testowanie

Program został przetestowany na różnych typach plików – graficznych, dźwiękowych, wykonawczych, tekstowych. W trakcie testowania (po kompresji i rozpakowaniu) nie zaobserwowano zmiany, ubytku żadnego bajta w pliku.

Największy plik dla którego przetestowano program ma 3,01 GB, kompresja pliku trwała ok 16 m, rozpakowywanie 40m 46s. Program wykazuje liniową złożoność czasową względem wielkości pliku. Plik o rozmiarze 1,5 GB rozpakowywał się około 20m, a o rozmiarze 0,3 GB 4m.

Maksymalny rozmiar pliku zależy od kompilatora (typ `int` może być realizowany jako zmienna dwu- lub czterobajtowa).

6. Wnioski

Program kompresji plików metodą kodowania Huffmana jest programem operującym na plikach binarnych. Pisanie go pozwoliło mi lepiej zrozumieć sposób strukturę plików. Nie jest ona skomplikowana jednak po raz pierwszy miałem większy projekt związany właśnie z ich obsługą i to właśnie to było najtrudniejsze w projekcie.

Źródła

Ron Rivest i Thomas H. Cormen „Wprowadzenie do algorytmów”

Huffman

4

Wygenerowano przez Doxygen 1.8.17

1 Indeks struktur danych	1
1.1 Struktury danych	1
2 Indeks plików	3
2.1 Lista plików	3
3 Dokumentacja struktur danych	5
3.1 Dokumentacja struktury min_heap	5
3.1.1 Opis szczegółowy	5
3.2 Dokumentacja struktury node	5
3.2.1 Opis szczegółowy	6
4 Dokumentacja plików	7
4.1 Dokumentacja pliku Main.c	7
4.2 Dokumentacja pliku sources.c	7
4.2.1 Dokumentacja funkcji	8
4.2.1.1 add_to_tree()	8
4.2.1.2 build_min_heap()	8
4.2.1.3 build_min_heap_i()	8
4.2.1.4 codematrix()	9
4.2.1.5 compress()	9
4.2.1.6 delete_tree()	10
4.2.1.7 extract_min()	10
4.2.1.8 heap_insert()	10
4.2.1.9 min_heapify()	10
4.2.1.10 read_file_coded()	11
4.2.1.11 read_file_uncoded()	11
4.2.1.12 readme()	12
4.2.1.13 save_to_file()	12
4.3 Dokumentacja pliku sources.h	12
4.3.1 Dokumentacja funkcji	13
4.3.1.1 add_to_tree()	13
4.3.1.2 build_min_heap()	13
4.3.1.3 build_min_heap_i()	14
4.3.1.4 codematrix()	14
4.3.1.5 compress()	15
4.3.1.6 delete_tree()	15
4.3.1.7 extract_min()	15
4.3.1.8 heap_insert()	16
4.3.1.9 min_heapify()	16
4.3.1.10 number_to_string()	16
4.3.1.11 read_file_coded()	17
4.3.1.12 read_file_uncoded()	17

4.3.1.13 readme()	17
4.3.1.14 save_to_file()	18

Indeks	21
---------------	-----------

Rozdział 1

Indeks struktur danych

1.1 Struktury danych

Tutaj znajdują się struktury danych wraz z ich krótkimi opisami:

min_heap	< Struktura kopca (statycznego)/ Kolejka priorytetowa typu min	5
node	< struktura węzła drzewa	5

Rozdział 2

Indeks plików

2.1 Lista plików

Tutaj znajduje się lista wszystkich udokumentowanych plików z ich krótkimi opisami:

Main.c	7
sources.c	7
sources.h	12

Rozdział 3

Dokumentacja struktur danych

3.1 Dokumentacja struktury min_heap

< Struktura kopca (statycznego)/ Kolejka priorytetowa typu min

```
#include <sources.h>
```

Pola danych

- **node * tab** [256]
tablica przechowująca wierzchołki drzewa
- **int lengh**
długość aktualnej kolejki

3.1.1 Opis szczegółowy

< Struktura kopca (statycznego)/ Kolejka priorytetowa typu min

Dokumentacja dla tej struktury została wygenerowana z pliku:

- **sources.h**

3.2 Dokumentacja struktury node

< struktura węzła drzewa

```
#include <sources.h>
```

Pola danych

- unsigned int **frequency**
częstość występowania bitu
- int **bajt**
bit którego dotyczy powyższa częstość lub -1 jeśli w węźle nie znajduje się żaden
- struct **node** * **left_n**
wskaźnik na lewego potomka
- struct **node** * **right_n**
wskaźnik na prawego potomka

3.2.1 Opis szczegółowy

< struktura węzła drzewa

Dokumentacja dla tej struktury została wygenerowana z pliku:

- **sources.h**

Rozdział 4

Dokumentacja plików

4.1 Dokumentacja pliku Main.c

```
#include <stdlib.h>
#include <stdio.h>
#include "sources.h"
```

Funkcje

- int **main** (int argc, char *argv[])

4.2 Dokumentacja pliku sources.c

```
#include <stdlib.h>
#include <stdio.h>
#include "sources.h"
```

Funkcje

- void **delete_tree** (node *root)
- void **min_heapify** (min_heap *A, int i)
- void **heap_insert** (min_heap *A, node *key)
- min_heap **build_min_heap** (node *A, int n)
- min_heap **build_min_heap_i** (node *A, int n)
- node * **extract_min** (min_heap *A)
- node * **add_to_tree** (node *leafs)
- node * **read_file_uncoded** (char *nazwa, node *tab)
- void **codematrix** (node *A, int **code, int i)
zwraca tablice(słownik) kodów huffmana na podstawie drzewa huffmna
- int **save_to_file** (char *iname, char *oname, int **codevector, node *leafs_huffman_tree)
- void **number_to_string** (char number, char string_n[1280], int *w)
- int **find_equivalent** (node *root, char code[1280], int *w)
- int **read_file_coded** (char *iname, char *oname)
- int **compress** (char *iname, char *oname)
- void **readme** ()

4.2.1 Dokumentacja funkcji

4.2.1.1 add_to_tree()

```
node* add_to_tree (
    node * leafs )
```

Buduje drzewo kodów Huffmna

Parametry

in, out	leafs	Tablica 256 liści drzewa
---------	-------	--------------------------

Zwraca

Wskaźnik na korzeń drzewa

4.2.1.2 build_min_heap()

```
min_heap build_min_heap (
    node * A,
    int n )
```

Buduje kopiec typu min z podanej elementowej tablicy

Parametry

in, out	A	Wskaźnik na kolejkę
in	n	ilość elementów tablicy

Zwraca

zwraca kolejkę typu min

4.2.1.3 build_min_heap_i()

```
min_heap build_min_heap_i (
    node * A,
    int n )
```

Buduje kopiec typu min z podanej elementowej tablicy

Parametry

in, out	<i>A</i>	Wskaźnik na kolejkę
in	<i>n</i>	ilość elementów tablicy

Zwraca

zwraca kolejkę typu min* zwraca kolejkę typu min

4.2.1.4 codematrix()

```
void codematrix (
    node * A,
    int ** code,
    int i )
```

zwraca tablice(słownik) kodów huffmana na podstawie drzewa huffmna

Funkcja tworzy słownik w którym dla i-tego wiersza występuje kod bajtu o wartości i, kod jest ciągiem 0 i 1 zakończonym liczbą 2

Parametry

in	<i>A</i>	Wskaźnik na korzeń drzewa modów Huffmana
in, out	<i>code</i>	Tablica 257x256 w której zwracany będzie słownik
in	<i>i</i>	Zmienna pomocnicza, domyślnie 0

4.2.1.5 compress()

```
int compress (
    char * iname,
    char * oname )
```

Funkcja tworzy obsługującą funkcjonalność kompresji pliku

Parametry

in	<i>iname</i>	Nazwa pliku wejściowego (niezakodowanego)
in	<i>oname</i>	Nazwa pliku wyjściowego

4.2.1.6 delete_tree()

```
void delete_tree (
    node * root )
```

Funkcja usuwa drzewo bez liści

Parametry

in, out	root	Wskaźnik na korzeń drzewa
---------	------	---------------------------

4.2.1.7 extract_min()

```
node* extract_min (
    min_heap * A )
```

Zciąga z kopca element o najmniejszej częstości występowania

Parametry

in, out	A	Wskaźnik na kolejkę
---------	---	---------------------

Zwraca

zwraca element o najmniejszej częstości występowania

4.2.1.8 heap_insert()

```
void heap_insert (
    min_heap * A,
    node * key )
```

Funkcja wstawia klucz do kopca

Parametry

in, out	A	Wskaźnik na kolejkę
in	key	Klucz który należy wstawić do kolejki

4.2.1.9 min_heapify()

```
void min_heapify (
```

```
    min_heap * A,  
    int i )
```

Funkcja przywraca własności kopca typu min

Parametry

in, out	A	Wskaźnik na kolejkę
in	i	element tablicy dla którego wywołujemy funkcję

4.2.1.10 read_file_coded()

```
int read_file_coded (  
    char * iname,  
    char * oname )
```

Funkcja tworzy obsługującą funkcjonalność dekompresji pliku

Parametry

in	iname	Nazwa pliku wejściowego (niezakodowanego)
in	oname	Nazwa pliku wyjściowego

Zwraca

0 w przypadku poprawnego wykonania funkcji

4.2.1.11 read_file_uncoded()

```
node* read_file_uncoded (  
    char * nazwa,  
    node * tab )
```

Funkcja uzupełnia dane w tablicy liści dotyczące częstości występowania bajtów na podswie pliku

Parametry

in, out	tab	Tablica 256 liści drzewa
in	nazwa	Nazwa pliku dla którego należy zrobić statystykę

Zwraca

Wskaźnik na przyjętą tablicę

4.2.1.12 readme()

```
void readme ( )
```

Funkcja tworzy plik readme

4.2.1.13 save_to_file()

```
int save_to_file (
    char * iname,
    char * oname,
    int ** codevector,
    node * leafs_huffman_tree )
```

Funkcja tworzy zakodowany plik razem z danymi potrzebnymi do jego odczytu

Parametry

in	<i>iname</i>	Nazwa pliku wejściowego (niezakodowanego)
in	<i>oname</i>	Nazwa pliku wyjściowego
in	<i>codevector</i>	słownik kodów Huffmana
in	<i>leafs_huffman_tree</i>	Tablica liści drzewa huffmana

Zwraca

0 w przypadku poprawnego wykonania funkcji

4.3 Dokumentacja pliku sources.h

Struktury danych

- struct **node**
< struktura węzła drzewa
- struct **min_heap**
< Struktura kopca (statycznego)/ Kolejka priorytetowa typu min

Definicje typów

- typedef struct **node** node
< struktura węzła drzewa
- typedef struct **min_heap** min_heap
< Struktura kopca (statycznego)/ Kolejka priorytetowa typu min

Funkcje

- void **delete_tree** (**node** *root)
- void **min_heapify** (**min_heap** *A, int i)
- void **heap_insert** (**min_heap** *A, **node** *key)
- **min_heap build_min_heap** (**node** *A, int n)
- **min_heap build_min_heap_i** (**node** *A, int n)
- **node** * **extract_min** (**min_heap** *A)
- **node** * **add_to_tree** (**node** *leafs)
- **node** * **read_file_uncoded** (char *nazwa, **node** *tab)
- void **codematrix** (**node** *A, int **code, int i)
zwraca tablice(słownik) kodów huffmana na podstawie drzewa huffmna
- void **number_to_string** (char number, char *string_n, int *w)
- void **readme** ()
- int **compress** (char *iname, char *oname)
- int **save_to_file** (char *iname, char *oname, int **codevector, **node** *leafs_huffman_tree)
- int **read_file_coded** (char *iname, char *oname)

4.3.1 Dokumentacja funkcji

4.3.1.1 add_to_tree()

```
node* add_to_tree (
    node * leafs )
```

Buduje drzewo kodów Huffmna

Parametry

in, out	leafs	Tablica 256 liści drzewa
---------	-------	--------------------------

Zwraca

Wskaźnik na korzeń drzewa

4.3.1.2 build_min_heap()

```
min_heap build_min_heap (
    node * A,
    int n )
```

Buduje kopiec typu min z podanej elementowej tablicy

Parametry

<code>in, out</code>	<code>A</code>	Wskaźnik na kolejkę
<code>in</code>	<code>n</code>	ilość elementów tablicy

Zwraca

zwraca kolejkę typu min

4.3.1.3 build_min_heap_i()

```
min_heap build_min_heap_i (
    node * A,
    int n )
```

Buduje kopiec typu min z podanej elementowej tablicy

Parametry

<code>in, out</code>	<code>A</code>	Wskaźnik na kolejkę
<code>in</code>	<code>n</code>	ilość elementów tablicy

Zwraca

zwraca kolejkę typu min* zwraca kolejkę typu min

4.3.1.4 codematrix()

```
void codematrix (
    node * A,
    int ** code,
    int i )
```

zwraca tablice(słownik) kodów huffmana na podstawie drzewa huffmna

Funkcja tworzy słownik w którym dla i-tego wiersza występuje kod bajtu o wartości i, kod jest ciągiem 0 i 1 zakończonym liczbą 2

Parametry

<code>in</code>	<code>A</code>	Wskaźnik na korzeń drzewa modów Huffmana
<code>in, out</code>	<code>code</code>	Tablica 257x256 w której zwracany będzie słownik
<code>in</code>	<code>i</code>	Zmienna pomocnicza, domyślnie 0

4.3.1.5 compress()

```
int compress (
    char * iname,
    char * oname )
```

Funkcja tworzy obsługującą funkcjonalność kompresji pliku

Parametry

in	<i>iname</i>	Nazwa pliku wejściowego (niezakodowanego)
in	<i>oname</i>	Nazwa pliku wyjściowego

4.3.1.6 delete_tree()

```
void delete_tree (
    node * root )
```

Funkcja usuwa drzewo bez liści

Parametry

in, out	<i>root</i>	Wskaźnik na korzeń drzewa
---------	-------------	---------------------------

4.3.1.7 extract_min()

```
node* extract_min (
    min_heap * A )
```

Zciąga z kopca element o najmniejszej częstości występowania

Parametry

in, out	<i>A</i>	Wskaźnik na kolejkę
---------	----------	---------------------

Zwraca

zwraca element o najmniejszej częstości występowania

4.3.1.8 heap_insert()

```
void heap_insert (
    min_heap * A,
    node * key )
```

Funkcja wstawia klucz do kopca

Parametry

in, out	<i>A</i>	Wskaźnik na kolejkę
in	<i>key</i>	Klucz który należy wstawić do kolejki

4.3.1.9 min_heapify()

```
void min_heapify (
    min_heap * A,
    int i )
```

Funkcja przywraca własności kopca typu min

Parametry

in, out	<i>A</i>	Wskaźnik na kolejkę
in	<i>i</i>	element tablicy dla którego wywołujemy funkcję

4.3.1.10 number_to_string()

```
void number_to_string (
    char number,
    char * string_n,
    int * w )
```

Pomocnicza funkcja zamiany liczb na ciąg znaków 0 i 1 zakończony 2 jako znakiem końca ciągu

Parametry

in	<i>number</i>	Liczba którą należy zmienić
in, out	<i>string</i> _↔ <i>_n</i>	Zwracany ciąg
in, out	<i>w</i>	Pierwszy niezapisany element ciągu

4.3.1.11 read_file_coded()

```
int read_file_coded (
    char * iname,
    char * oname )
```

Funkcja tworzy obsługującą funkcjonalność dekompresji pliku

Parametry

in	<i>iname</i>	Nazwa pliku wejściowego (niezakodowanego)
in	<i>oname</i>	Nazwa pliku wyjściowego

Zwraca

0 w przypadku poprawnego wykonania funkcji

4.3.1.12 read_file_uncoded()

```
node* read_file_uncoded (
    char * nazwa,
    node * tab )
```

Funkcja uzupełnia dane w tablicy liści dotyczące częstości występowania bajtów na podstawie pliku

Parametry

in, out	<i>tab</i>	Tablica 256 liści drzewa
in	<i>nazwa</i>	Nazwa pliku dla którego należy zrobić statystykę

Zwraca

Wskaźnik na przyjętą tablicę

4.3.1.13 readme()

```
void readme ( )
```

Funkcja tworzy plik readme

4.3.1.14 `save_to_file()`

```
int save_to_file (
    char * iname,
    char * oname,
    int ** codevector,
    node * leafs_huffman_tree )
```

Funkcja tworzy zakodowany plik razem z danymi potrzebnymi do jego odczytu

Parametry

in	<i>iname</i>	Nazwa pliku wejściowego (niezakodowanego)
in	<i>oname</i>	Nazwa pliku wyjściowego
in	<i>codevector</i>	słownik kodów Huffmana
in	<i>leafs_huffman_tree</i>	Tablica liści drzewa huffmana

Zwraca

0 w przypadku poprawnego wykonania funkcji

Indeks

add_to_tree
 sources.c, 8
 sources.h, 13

build_min_heap
 sources.c, 8
 sources.h, 13

build_min_heap_i
 sources.c, 8
 sources.h, 14

codematrix
 sources.c, 9
 sources.h, 14

compress
 sources.c, 9
 sources.h, 15

delete_tree
 sources.c, 9
 sources.h, 15

extract_min
 sources.c, 10
 sources.h, 15

heap_insert
 sources.c, 10
 sources.h, 15

Main.c, 7
min_heap, 5
min_heapify
 sources.c, 10
 sources.h, 16

node, 5
number_to_string
 sources.h, 16

read_file_coded
 sources.c, 11
 sources.h, 16

read_file_uncoded
 sources.c, 11
 sources.h, 17

readme
 sources.c, 11
 sources.h, 17

save_to_file

 sources.c, 12
 sources.h, 17
sources.c, 7
 add_to_tree, 8
 build_min_heap, 8
 build_min_heap_i, 8
 codematrix, 9
 compress, 9
 delete_tree, 9
 extract_min, 10
 heap_insert, 10
 min_heapify, 10
 read_file_coded, 11
 read_file_uncoded, 11
 readme, 11
 save_to_file, 12
sources.h, 12
 add_to_tree, 13
 build_min_heap, 13
 build_min_heap_i, 14
 codematrix, 14
 compress, 15
 delete_tree, 15
 extract_min, 15
 heap_insert, 15
 min_heapify, 16
 number_to_string, 16
 read_file_coded, 16
 read_file_uncoded, 17
 readme, 17
 save_to_file, 17