CS213/293 Data Structure and Algorithms 2024

Lecture 11: Data compression

Instructor: Ashutosh Gupta

IITB India

Compile date: 2024-10-13



You must have used Zip, which reduces the space a file uses.

How does Zip work?



Fixed-length vs. Variable-length encoding

- Fixed-length encoding. Example: An 8-bit ASCII code encodes each character in a text file.
- ▶ Variable-length encoding: each character is given a different bit length encoding.
- ▶ We may save space by assigning fewer bits to the characters that occur more often.
- We may have to assign some characters more than 8-bit representation.



Example: Variable-length encoding

Example 11.1

Consider text: "agra"

- In a text file, the text will take 32 bits of space.
 - 01100001011001110111001001100001
- ▶ There are only three characters. Let us use encoding, a = "0", g = "10", and r = "11". The text needs six bits.
 - 010110

Exercise 11.1

Are the six bits sufficient?

Commentary: If the encoding depends on the text content, we also need to record the encoding along with the text.

Example: decoding variable-length encoding

Example 11.2

Consider encoding a = "0", g = "10", and r = "11" and the following encoding of a text.

101100001110

The text is "graaaarg".

We scan the encoding from the left. As soon as a match is found, we start matching the next symbol.



Example: decoding bad variable-length encoding

Example 11.3

Consider encoding a = "0", g = "01", and r = "11" and the following encoding of a text.

0111000011001

We cannot tell if the text starts with a "g" or an "a".

Prefix condition: Encoding of a character cannot be a prefix of encoding of another character.



Encoding trie

Definition 11.1

An encoding trie is a binary trie that has the following properties.

- Each terminating leaf is labeled with an encoded character.
- The left child of a node is labeled 0 and the right child of a node is labeled 1

Exercise 11.2

Show: An encoding trie ensures that the prefix condition is not violated.



Character encoding/codewords: C = 00, A = 010, R = 011, D = 10, and B = 11. Example: Decoding from a Trie



Encoding: 01011011010000101001011011010

Text: ABRACADABRA

Encoding length

Example 11.4

Let us encode ABRACADABRA using the following two tries.



Encoding:(29 bits) 01011011010 0001010 01011011010 Encoding:(24 bits) 00111000 01000011 00111000



В

Drawing with tries without labels

Since we know the label of an internal node by observing that a node is a left or right child, we will not write the labels.



Commentary: We can assign any bit to a node as long as the sibling will use a different bit.



Topic 11.1

Optimal compression

Different tries will result in different compression levels.

Design principle: We encode a character that occurs more often with fewer bits.



frequency

Definition 11.2

The frequency f_c of a character c in a text T is the number of times c occurs in T.

Example 11.5

The frequencies of the characters in ABRACADABRA are as follows.

►
$$f_A = 5$$

•
$$f_B = 2$$

•
$$f_R = 2$$

•
$$f_C = 1$$

•
$$f_D = 1$$



Characters encoding length

Definition 11.3 The encoding length l_c of a character c in a trie is the number of bits needed to encode c.

Example 11.6



In the left trie, the encoding length of the characters are as follows.

Weighted path length == number of encoded bits

The total number of bits needed to store a text is

Example 11.7

The number of bits needed for ABRACADABRA using the left trie is the following sum.

$$f_A * I_A + f_C * I_C + f_D * I_D + f_R * I_R + f_B * I_B$$

$$= 5 * 2 + 1 * 3 + 1 * 3 + 2 * 2 + 2 * 2 = 24$$

Is this the best trie for compression? How can we find the best trie?

Huffman encoding

Algorithm 11.1: HUFFMAN (Integers f_{c_1}, \ldots, f_{c_k}) 1 for $i \in [1, k]$ do CREATENODE(Value, LeftChild, RightChild) $N := CREATENODE(c_i, Null, Null); <$ 2 is a constructor of a node. 3 $T_i := \text{CREATENODE}(f_{c_i}, N, Null);$ 4 return BuildTree $(T_1, ..., T_k)$ Exercise 11.3 a. Is BuildTree tail recursive? Algorithm 11.2: BUILDTREE(Nodes $T_1, ..., T_k$) b. How should we resolve 1 if k == 1 then non-determinism if there is a tie return T_1 2 in finding the minimum? **3** Find T_i and T_i such that $value(T_i)$ and $value(T_i)$ are minimum:

- 4 $T_i := CREATENODE(value(T_i) + value(T_j), T_i, T_j);$
- **5 return** $BuildTree(T_1, ..., T_{j-1}, T_{j+1}, ..., T_k)$

Running time analysis of Huffman encoding

We need to find minimums repeatedly. We use a heap to store the values of the roots.

Running time analysis

- BuildTree will be recursively called k times.
- In each recursive call, we need to call
 - two deleteMins for removing two trees and
 - an insertion for the new tree

in the heap.

Total running time

$$\sum_{i=k}^1 O(\log i) = O(k \log k)$$

Commentary: We have proven the above equality in tutorial problems!

Example: Huffman encoding

- Example 11.8
- After initialization.

We choose nodes labeled with 1 to join and create a larger tree.

Example: Huffman encoding(2)

After the next recursive step

After another recursive step:

Example: Huffman encoding(3)

After the final recursive step:

We scrub the frequency labels.

Exercise 11.4

How many bits do we need to encode ABRACADABRA?

Using Huffman in compression

- To compress a file, we need to compute the frequencies of the symbols.
- The number of symbols may be constant with respect to the file size.
- Therefore, the cost of computing Huffman is constant time if the frequencies are given.

Frequency (alternative definition!)

Definition 11.4 (Equivalent definition)

The frequency f_c of a character c in a text T is the fraction (or %) of times c occurs in T.

Example 11.9

The frequencies of the characters in ABRACADABRA are as follows.

- $f_A = 5/11$
- ▶ $f_B = 2/11$
- ▶ $f_R = 2/11$
- ▶ $f_C = 1/11$
- $f_D = 1/11$

Huffman can work with the fractions without any change. $_{\mbox{(Why?)}}$

Topic 11.2

Proof of optimality of Huffman encoding

Is Huffman optimal?

Exercise 11.5

Let us suppose a file contains a, b, c, and d with frequencies 25%, 25%, 25%, and 25% respectively.

- a. Should you be able to compress this file?
- b. Do Huffman codes compress this file?

We need to prove that Huffman encoding indeed produces optimal encoding.

Minimum weighted path length

Definition 11.5

Given frequencies $f_{c_1}, ..., f_{c_k}$, minimum weighted path length $MWPL(f_{c_1}, ..., f_{c_k})$ is the weighted path length for the encoding trie that encodes $c_1, ..., c_k$ and the sum is minimum.

We say a trie is a witness of $MWPL(f_{c_1}, ..., f_{c_k})$ if it encodes $c_1, ..., c_k$ and it produces encoding of length $MWPL(f_{c_1}, ..., f_{c_k})$ for a text with frequencies $f_{c_1}, ..., f_{c_k}$

Example 11.10

We have seen MWPL(5, 2, 2, 1, 1) = 23.

The witness trie is on the left.

Commentary: The definition of MWPL does not mention the trie. It is the property of the frequency distribution.

A recursive relation

Theorem 11.1 $MWPL(f_{c_1}, ..., f_{c_k}) \leq f_{c_1} + f_{c_2} + MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k})$

Proof.

Let trie T be a witness of $MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k})$ containing a node labeled with $f_{c_1} + f_{c_2}$ with a terminal child.

A recursive relation(2)

Proof(contd.)

We construct a trie for frequencies $f_{c_1}, ..., f_{c_k}$ such that the weighted path length of the trie is $f_{c_1} + f_{c_2} + MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k})$.

Therefore, $MWPL(f_{c_1}, ..., f_{c_k})$ must be less than equal to the above expression.

27

Example: $MWPL(f_{c_1}, ..., f_{c_k}) \leq f_{c_1} + f_{c_2} + MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k})$ Example 11.11

Witness for MWPL(5, 2, 2, 2)

The weighted path length of the above is 1 + 4 + MWPL(5, 2, 2, 2)

The witness of MWPL(1, 4, 2, 2, 2) must have weighted path length \leq the above right trie.

 $MWPL(1, 4, 2, 2, 2) \le 1 + 4 + MWPL(5, 2, 2, 2)$

Reverse recursive relation

Theorem 11.2

If f_{c_1} and f_{c_2} are the minimum two, $MWPL(f_{c_1}, ..., f_{c_k}) = f_{c_1} + f_{c_2} + MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k})$. Proof.

There is a witness of $MWPL(f_{c_1}, ..., f_{c_k})$ where the parents of c_1 and c_2 are siblings. (Why?)

Commentary: Explaining why: Show that the smallest frequency symbol can always be moved to the last level to improve weighted path length. Furthermore, since there must be a sibling at the last level, the second last frequency symbol can also be moved to the sibling to improve the weighted path length.

\sim	0	0	\sim
(cc)	7 a 1	~	()
	чт ,	1 24	
~	~	~	\sim

Instructor: Ashutosh Gupta

Reverse recursive relation(2)

Proof(contd.)

We construct a tree for frequencies $f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k}$ such that the weighted path length of the tree is $MWPL(f_{c_1}, ..., f_{c_k}) - f_{c_1} - f_{c_2}$.

Therefore, $MWPL(f_{c_1}, ..., f_{c_k}) - f_{c_1} - f_{c_2} \ge MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k}).$

Due to the previous theorem, $MWPL(f_{c_1}, ..., f_{c_k}) = f_{c_1} + f_{c_2} + MWPL(f_{c_1} + f_{c_2}, f_{c_3}, ..., f_{c_k})$.

Witness for MWPL(5,2,2,2). Since 2 and 2 are the least two frequencies, they are on the longest path.

The weighted path length of the above is MWPL(5, 2, 2, 2) - 2 - 2

The witness of MWPL(5, 2, 4) must have weighted path length \leq the above right trie.

 $MWPL(5,2,4) \leq MWPL(5,2,2,2) - 2 - 2$

Instructor: Ashutosh Gupta

Proof compatible $\operatorname{BuildTree}$

Our BUILDTREE does not follow the pattern of updates of theorem 11.2. So, writing a proof over the algorithm is hard. We will consider the following version of BUILDTREE when writing the proof.

Algorithm 11.3: BUILDTREE2(Nodes $T_1, ..., T_k$)

- 1 if k == 1 then
- 2 return T_1
- 3 Find T_i and T_j such that $value(T_i)$ and $value(T_j)$ are minimum;
- 4 $T_{new} := CREATENODE(value(T_i) + value(T_j), Null, Null);$
- 5 $T := BUILDTREE2(T_1, ..., T_{i-1}, T_{i+1}, ..., T_{j-1}, T_{j+1}, ..., T_k, T_{new});$
- 6 $left(T_{new}) := T_i;$
- 7 $right(T_{new}) := T_j;$
- 8 return T

Exercise 11.6

Show that algorithms 11.2 and 11.3 are equivalent.

Correctness of $\operatorname{Huffman}$

Theorem 11.3

HUFFMAN $(f_{c_1}, ..., f_{c_k})$ always returns a tree that is a witness of $MWPL(f_{c_1}, ..., f_{c_k})$.

Proof.

We assume HUFFMAN is calling BUILDTREE2. We prove inductively over k in the call of BUILDTREE2($T_1, ..., T_k$).

Base case:

Trivial. There is a single tree with a single node and we return the node.

Induction step:

We assume the recursive call BUILDTREE2 returns witness T of $MWPL(value(T_1), ..., value(T_{i-1}), value(T_{i+1}), ..., value(T_{j-1}), value(T_{j+1}), ..., value(T_k), value(T_{new}))$.

Therefore, T is a witness of $MWPL(value(T_1), ..., value(T_{i-1}), value(T_{i-1}), ..., value(T_{j-1}), value(T_{j+1}), ..., value(T_k), value(T_i) + value(T_j)).$

Subsequently, we insert nodes T_i and T_j in T according to the scheme of theorem 11.2.

Therefore, the T at line 10 is a witness of $MWPL(value(T_1), ..., value(T_k))$.

33

Topic 11.3

Repeated string (LZ77)

LZ77 for repeated string

In LZ77, if a string is repeated within the sliding window on the input stream, the repeated occurrence is replaced by a reference, which is a pair of the offset and length of the string.

The references are viewed as yet another symbol on the input stream.

Example 11.13

Before encoding ABRACADABRA using a trie, the string will be transformed to

ABRACAD[4,7].

We run Huffman on the above string.

Multiple repetitions

Example 11.14

Consider the following input text of 16 characters.

abababababababab

We will transform the text as follows.

ab[14, 2]

Topic 11.4

DEFLATE

Practical Huffman

When we compress a file, we do not compute the frequencies for the entire file in one go.

- We compute the encoding trie of a block of bytes.
- ▶ We check if the data allows compression, if it does not we do not compress the block
- If the block is small, we use a precomputed encoding trie.

Exercise 11.7

How many bits are needed per character for 8 characters if frequencies are all equal?

DEFLATE

The Linux utility gzip uses the DEFLATE algorithm for compression, which combines Huffman encoding and the LZ77 algorithm.

DEFLATE compresses one file in blocks. Each block may be compressed in one of three modes.

- No compression
- Dynamically computed Huffman coding
- Fixed encoding

To compress multiple files, we first use a tar utility that concatenates the files into one file.

Let us look the content of the gzip file

Example 11.15

Let us consider a file "name.txt" that contains "abracadabra".

We compress the file using the following command.

gzip -kf name.txt

The command will generate file name.txt.gz. We may view the content of the file as follows.

xxd -b name.txt.gz

The contents are displayed in the next slide.

uncompressed filesize----

A printing problem

The print in the previous slide is generated via the following command.

```
xxd -b name.txt.gz
```

It prints the bytes from MSB to LSB. So, we are seeing each byte reversed.

We need to reverse each byte to see the actual DEFLATE stream.

Example 11.16

The first byte of the DEFLATE stream is printed as 01001011, which should be read as 11010010.

The following is the DEFLATE stream of 12 bytes from the previous slide.

b. How does gzip identify repeated patterns?

Topic 11.5

Tutorial problems

Single-bit Huffman code

Exercise 11.9

a. In a Huffman code instance, show that if there is a character with a frequency greater than $\frac{2}{5}$ then there is a codeword of length 1.

b. Show that if all frequencies are less than $\frac{1}{3}$ then there is no codeword of length 1.

Predictable text

Exercise 11.10

Suppose that there is a source that has three characters a,b,c. The output of the source cycles in the order of a,b,c followed by a again, and so on. In other words, if the last output was a b, then the next output will either be a b or a c. Each letter is equally probable. Is the Huffman code the best possible encoding? Are there any other possibilities? What would be the pros and cons of this?

Compute Huffman code tree

Exercise 11.11

Given the following frequencies, compute the Huffman code tree.

а	20
d	7
g	8
j	4
b	6
е	25
h	8
k	2
с	6
f	1
i	12
Ι	1

Topic 11.6

Problems

Exercise: Huffman tree for fixed encoding of DEFLATE

Example 11.17

Draw the Huffman tree for the fixed encoding used in DEFLATE.

In DEFLATE, there are 288 symbols. 0-255 are the input byte, 257-287 are the length symbols, and the 256th symbol is for the end of a block. The following is from DEFLATE RFC.

The Huffman codes for the two alphabets are fixed, and are not represented explicitly in the data. The Huffman code lengths for the literal/length alphabet are:

Lit	Value	Bits	Codes
0	- 143	8	00110000 through 10111111
144	- 255	9	110010000 through 11111111
256	- 279	7	0000000 through 0010111
280	- 287	8	11000000 through 11000111

End of Lecture 11

