Lecture No.19 Huffman Encoding

CC-213 Data Structures Department of Computer Science University of the Punjab

Slides modified very slightly from the late Dr. Sohail Aslam's lectures at VU

Other Uses of Binary Trees

- Data compression plays a significant role in computer networks.
- To transmit data to its destination faster, it is necessary to either increase the data rate of the transmission media or to simply send less data.
- Improvements with regard to the transmission media has led to increase in the rate.
- The other options is to send less data by means of data compression.
- Compression methods are used for text, images, voice and other types of data (space probes).

- Huffman code is method for the compression for standard text documents.
- It makes use of a binary tree to develop codes of varying lengths for the letters used in the original message.
- Huffman code is also part of the JPEG image compression scheme.
- The algorithm was introduced by David Huffman in 1952 as part of a course assignment at MIT.

- To understand Huffman encoding, it is best to use a simple example.
- Encoding the 32-character phrase: "traversing threaded binary trees",
- If we send the phrase as a message in a network using standard 8-bit ASCII codes, we would have to send 8*32= 256 bits.
- Using the Huffman algorithm, we can send the message with only 116 bits.

- List all the letters used, including the "space" character, along with the frequency with which they occur in the message.
- Consider each of these (character,frequency) pairs to be nodes; they are actually leaf nodes, as we will see.
- Pick the two nodes with the lowest frequency, and if there is a tie, pick randomly amongst those with equal frequencies.

- Make a new node out of these two, and make the two nodes its children.
- This new node is assigned the sum of the frequencies of its children.
- Continue the process of combining the two nodes of lowest frequency until only one node, the root, remains.

Original text: *traversing threaded binary trees* size: 33 characters (space and newline)

| NL: 1 | i: | 2 |
|-------|------------|---|
| SP: 3 | n : | 2 |
| a: 3 | r: | 5 |
| b: 1 | S : | 2 |
| d: 2 | t : | 3 |
| e: 5 | V : | 1 |
| g: 1 | V : | 1 |
| | | |



There a number of ways to combine nodes. We have chosen just one such way.















- Huffman code is method for the compression for standard text documents.
- It makes use of a binary tree to develop codes of varying lengths for the letters used in the original message.
- Huffman code is also part of the JPEG image compression scheme.
- The algorithm was introduced by David Huffman in 1952 as part of a course assignment at MIT.

- To understand Huffman encoding, it is best to use a simple example.
- Encoding the 32-character phrase: "traversing threaded binary trees",
- If we send the phrase as a message in a network using standard 8-bit ASCII codes, we would have to send 8*32= 256 bits.
- Using the Huffman algorithm, we can send the message with only 116 bits.

- List all the letters used, including the "space" character, along with the frequency with which they occur in the message.
- Consider each of these (character,frequency) pairs to be nodes; they are actually leaf nodes, as we will see.
- Pick the two nodes with the lowest frequency, and if there is a tie, pick randomly amongst those with equal frequencies.

- Make a new node out of these two, and make the two nodes its children.
- This new node is assigned the sum of the frequencies of its children.
- Continue the process of combining the two nodes of lowest frequency until only one node, the root, remains.

Original text: *traversing threaded binary trees* size: 33 characters (space and newline)

| NL: 1 | i: 2 | |
|-------|------|--|
| SP: 3 | n: 2 | |
| a: 3 | r: 5 | |
| b: 1 | s: 2 | |
| d: 2 | t: 3 | |
| e: 5 | v: 1 | |
| g: 1 | v: 1 | |
| | | |



There a number of ways to combine nodes. We have chosen just one such way.















- List all the letters used, including the "space" character, along with the frequency with which they occur in the message.
- Consider each of these (character,frequency) pairs to be nodes; they are actually leaf nodes, as we will see.
- Pick the two nodes with the lowest frequency, and if there is a tie, pick randomly amongst those with equal frequencies.

- Make a new node out of these two, and make the two nodes its children.
- This new node is assigned the sum of the frequencies of its children.
- Continue the process of combining the two nodes of lowest frequency until only one node, the root, remains.

- Start at the root. Assign 0 to left branch and 1 to the right branch.
- Repeat the process down the left and right subtrees.
- To get the code for a character, traverse the tree from the root to the character leaf node and read off the 0 and 1 along the path.









Huffman character codes

| NL | | 10000 |
|-----|---------------|-------|
| SP | \Rightarrow | 1111 |
| а | \Rightarrow | 000 |
| b | \Rightarrow | 10001 |
| d | \Rightarrow | 0100 |
| е | \Rightarrow | 101 |
| g | \Rightarrow | 10010 |
| h | \Rightarrow | 10011 |
| i – | \Rightarrow | 0101 |
| n | \Rightarrow | 0110 |
| r | \Rightarrow | 110 |
| S | \Rightarrow | 0111 |
| t | \Rightarrow | 001 |
| V | | 11100 |
| у | | 11101 |

- Notice that the code is variable length.
- Letters with higher frequencies have shorter codes.
- The tree could have been built in a number of ways; each would yielded different codes but the code would still be minimal.

Original: traversing threaded binary trees

Encoded:

Original: *traversing threaded binary trees* With 8 bits per character, length is 264.

Compressed into 122 bits, 54% reduction.