

Huffman Coding Algorithm: Complete Implementation and Analysis

Computer Science Project
Adam Hajjaji - Zakariae Jali

December 12, 2025

Abstract

This report presents a complete implementation of the Huffman coding algorithm, a lossless data compression technique that assigns variable-length binary codes to characters based on their frequencies. The project includes both the core algorithmic implementation and an interactive JavaFX visualization tool. We detail the main concepts, data structures, implementation specifics, and provide the essential proofs of the algorithm's properties.

Contents

1	Introduction	3
1.1	Project Overview	3
1.2	Report Structure	3
2	Basic Concepts	3
2.1	Tree Properties	3
2.2	Prefix-Free Codes	3
3	Huffman Algorithm Description	3
3.1	Problem	3
3.2	Algorithm Steps	4
4	Implementation Architecture	4
4.1	Main Classes	4
4.1.1	Model Classes	4
4.1.2	Algorithm Classes	4
4.1.3	UI Class	4
5	Core Algorithm Implementation	4
5.1	Frequency Analysis	4
5.2	Tree Construction	5
5.3	Code Generation	5

6	Visualization Interface	6
6.1	JavaFX Application	6
6.2	Tree Drawing	6
7	Formal Proofs	6
7.1	Tree Properties	6
7.2	Huffman Algorithm Optimality	7
8	Implementation Details	7
8.1	Important Design Choices	7
8.1.1	Sorting Method	7
8.1.2	Special Character Handling	8
9	Results and Performance	8
9.1	Compression Example	8
9.2	Visualization Features	8
10	Conclusion	8
10.1	Project Success	8

1 Introduction

1.1 Project Overview

Huffman coding, developed by David A. Huffman in 1952, is a compression algorithm that assigns shorter binary codes to more frequent characters and longer codes to less frequent ones. This project implements the complete Huffman coding pipeline including frequency analysis, tree construction, encoding, decoding, and an educational visualization interface.

1.2 Report Structure

This report is organized as follows:

- **Section 2: Basic Concepts** - Essential theoretical background
- **Section 3: Algorithm Description** - Explanation of Huffman's algorithm
- **Section 4: Implementation Architecture** - Class structure
- **Section 5: Core Algorithm Implementation** - Code details
- **Section 6: Visualization Interface** - JavaFX GUI
- **Section 7: Formal Proofs** - Important mathematical proofs
- **Section 8: Conclusion**

2 Basic Concepts

2.1 Tree Properties

A **tree** is a connected graph without cycles. For Huffman coding, we use binary trees where:

- Each node has at most two children (left and right)
- Leaves represent characters
- The path from root to leaf gives the binary code
- Left edges = 0, right edges = 1

2.2 Prefix-Free Codes

A **prefix-free code** means no code is the beginning of another code. This allows decoding without special markers. Huffman codes are always prefix-free.

3 Huffman Algorithm Description

3.1 Problem

Given characters with frequencies, find binary codes that minimize:

$$\text{Total bits} = \sum (\text{frequency} \times \text{code length})$$

3.2 Algorithm Steps

1. Count frequency of each character
2. Create a leaf node for each character with its frequency
3. Repeatedly combine the two nodes with smallest frequencies
4. Make a new node with frequency = sum of the two
5. Continue until only one node remains (the root)
6. Assign codes: left = 0, right = 1

4 Implementation Architecture

4.1 Main Classes

The project has three main parts:

4.1.1 Model Classes

- `HuffmanNode`: Tree node with character and frequency
- `HuffmanTree`: Complete tree structure

4.1.2 Algorithm Classes

- `FrequencyCounter`: Counts character frequencies
- `TreeBuilder`: Builds the Huffman tree
- `CodeGenerator`: Creates binary codes

4.1.3 UI Class

- `HuffmanTreeVisualizer`: Interactive visualization tool

5 Core Algorithm Implementation

5.1 Frequency Analysis

```
1 public static HashMap<Character, Integer> getFrequency(List<Character>
2   text){
3     HashMap<Character, Integer> res = new HashMap<>();
4     for (int i=0; i<text.size(); i++){
5       Character curr_char = text.get(i);
6       if (res.containsKey(curr_char)){
7         res.put(curr_char, res.get(curr_char)+1);
8       }
9       else{
10         res.put(curr_char, 1);
11     }
12   }
13 }
```

```

11     }
12     return res;
13 }
```

Listing 1: Frequency Counter

5.2 Tree Construction

```

1 public static HuffmanTree buildHuffmanTree(String text){
2     HuffmanTree res = new HuffmanTree();
3     List<HuffmanNode> list_nodes = SortedNodesCreator.sort(text);
4     LinkedList<HuffmanNode> ll_nodes = new LinkedList<>(list_nodes);
5
6     while (ll_nodes.size()>1){
7         List<HuffmanNode> actualStep = new ArrayList<>(ll_nodes);
8         res.addConstructionSteps(actualStep);
9
10        HuffmanNode node1 = ll_nodes.removeFirst();
11        HuffmanNode node2 = ll_nodes.removeFirst();
12
13        int new_frequency = node1.getFrequency()+node2.getFrequency();
14        HuffmanNode new_node = new HuffmanNode(null,new_frequency);
15        new_node.setLeft(node1);
16        new_node.setRight(node2);
17
18        // Insert in correct position to keep sorted
19        ListIterator<HuffmanNode> it = ll_nodes.listIterator();
20        boolean found = false;
21        while (it.hasNext() && !found){
22            HuffmanNode compareNode = it.next();
23            if (compareNode.getFrequency() > new_frequency){
24                it.previous();
25                it.add(new_node);
26                found = true;
27            }
28        }
29        if (!found){
30            it.add(new_node);
31        }
32    }
33    res.setRoot(ll_nodes.get(0));
34    return res;
35 }
```

Listing 2: Tree Building

5.3 Code Generation

```

1 public static void generateCode(HuffmanNode tree,
2                                 HashMap<Character, String> codeMap,
3                                 String encoding){
4     if (tree.isLeaf()){
5         codeMap.put(tree.getCharacter(),encoding);
6         return;
7     }
8     generateCode(tree.getLeft(), codeMap, encoding+"0");
9 }
```

```

9     generateCode(tree.getRight(), codeMap, encoding+"1");
10 }

```

Listing 3: Code Generation

6 Visualization Interface

6.1 JavaFX Application

The visualization tool provides:

- File loading from text files
- Step-by-step tree construction animation
- Character path highlighting
- Compression statistics display
- Auto-play for demonstration

6.2 Tree Drawing

```

1 private void calculateTreePositions(HuffmanNode node,
2                                     Map<HuffmanNode, NodePosition>
3                                     positions,
4                                     double x, double y,
5                                     double offset, double
6                                     verticalSpacing){
7     if (node == null) return;
8     positions.put(node, new NodePosition(x, y));
9
10    if (node.getLeft() != null){
11        double childOffset = offset * 0.6;
12        calculateTreePositions(node.getLeft(), positions,
13                               x - childOffset, y + verticalSpacing,
14                               childOffset, verticalSpacing);
15    }
16    if (node.getRight() != null){
17        double childOffset = offset * 0.6;
18        calculateTreePositions(node.getRight(), positions,
19                               x + childOffset, y + verticalSpacing,
20                               childOffset, verticalSpacing);
21    }
22 }

```

Listing 4: Tree Position Calculation

7 Formal Proofs

7.1 Tree Properties

Trees have at least two leaves. Take the longest path in the tree. The endpoints of this path cannot connect to any other vertices (would make a longer path or a cycle). So they have degree 1, making them leaves. \square

[Tree Characterizations] For a graph with n vertices, these are equivalent:

1. Connected with no cycles
2. Connected with $n - 1$ edges
3. Has $n - 1$ edges and no cycles
4. Exactly one path between any two vertices

Proof. We show the main ideas:

(1) (2): Removing a leaf gives smaller tree. By induction, smaller tree has $(n - 1) - 1$ edges. Add back leaf gives $n - 1$ edges.

(2) (3): Connected graph with $n - 1$ edges can't have cycles (would have fewer edges).

(3) (1): Graph with $n - 1$ edges and no cycles must be connected.

(1) (4): In tree, exactly one path. If exactly one path between vertices, graph is connected and has no cycles. \square

7.2 Huffman Algorithm Optimality

Huffman's algorithm produces the optimal prefix-free code.

Proof. By induction on number of characters.

Base: With 2 characters, codes "0" and "1" are optimal.

Induction: Assume optimal for $n - 1$ characters. For n characters:

1. In any optimal tree, the two least frequent characters are siblings at deepest level
2. Combine them into one node with combined frequency
3. By induction, Huffman gives optimal tree for $n - 1$ characters
4. Splitting the combined node gives optimal tree for n characters

\square

8 Implementation Details

8.1 Important Design Choices

8.1.1 Sorting Method

We use merge sort to sort nodes by frequency:

```
1 public static List<HuffmanNode> mergeSort(List<HuffmanNode> A){  
2     if (A.size() <= 1){return A;}  
3     int mid = A.size()/2;  
4     return merge(mergeSort(A.subList(0, mid)),  
5                  mergeSort(A.subList(mid, A.size())));  
6 }
```

8.1.2 Special Character Handling

```
1 private String formatCharacter(char c){
2     switch (c){
3         case ' ': return " ";
4         case '\n': return " ";
5         case '\t': return " ";
6         case '\r': return " ";
7         default:
8             if (c < 32 || c > 126) return " ";
9             return String.valueOf(c);
10    }
11 }
```

9 Results and Performance

9.1 Compression Example

For the default text:

- Original: 305 characters \times 8 bits = 2440 bits
- Compressed: 1344 bits
- Savings: 44.9%
- Average code length: 4.41 bits per character

9.2 Visualization Features

The tool successfully shows:

1. Initial sorted nodes
2. Step-by-step merging
3. Final tree structure
4. Character encoding paths

10 Conclusion

10.1 Project Success

This implementation:

- Correctly implements Huffman coding
- Provides educational visualization
- Handles various text inputs
- Shows algorithm in action