

Алгоритм относится к группе “жадных” (greedy) алгоритмов и использует только частоту появления одинаковых символов во входном блоке данных. Ставит в соответствие символам входного потока, которые встречаются чаще, цепочку битов меньшей длины. И напротив, встречающимся редко - цепочку большей длины. Для сбора статистики требует двух проходов по входному блоку (также существуют однократные адаптивные варианты алгоритма).

Характеристики алгоритма Хаффмана [1]:

- Степени сжатия: 8, 1.5, 1 (лучшая, средняя, худшая степени)
- Симметричность по времени: 2:1 (за счёт того, что требует двух проходов по массиву сжимаемых данных)
- Характерные особенности: один из немногих алгоритмов, который не увеличивает размера исходных данных в худшем случае (если не считать необходимости хранить таблицу перекодировки вместе с файлом).

Основные этапы алгоритма сжатия с помощью кодов Хаффмана

1. Сбор статистической информации для последующего построения таблиц кодов переменной длины
 2. Построение кодов переменной длины на основании собранной статистической информации
 3. Кодирование (сжатие) данных с использованием построенных кодов
- Описанный выше алгоритм сжатия требует хранения и передачи вместе с кодированными данными дополнительной информации, которая позволяет однозначно восстановить таблицу соответствия кодируемых символов и кодирующих битовых цепочек.

Следует отметить, что в некоторых случаях можно использовать постоянную таблицу (или набор таблиц), которые заранее известны как при кодировании, так и при декодировании. Или же строить таблицу адаптивно в процессе сжатия и восстановления. В этих случаях хранение и передача дополнительной информации не требуется, а также отпадает необходимость в предварительном сборе статистической информации (этап 1).

В дальнейшем мы будем рассматривать только двухпроходную схему с явной передачей дополнительной информации о таблице соответствия кодируемых символов и кодирующих битовых цепочек.

Первый проход по данным: сбор статистической информации

Будем считать, что входные символы - это байты. Тогда сбор статистической информации будет заключаться в подсчёте числа появлений одинаковых байтов во входном блоке данных.

Построение кодов переменной длины на основании собранной статистической информации

Алгоритм Хаффмана основывается на создании двоичного дерева [3].

Изначально все узлы считаются листьями (конечными узлами), которые представляют символ (в нашем случае байт) и число его появлений. Для построения двоичного дерева используется очередь с приоритетами, где узлу с наименьшей частотой присваивается наивысший приоритет.

Алгоритм будет включать в себя следующие шаги:

Шаг 1. Помещаем все листья с ненулевым числом появлений в очередь с приоритетами (упорядочиваем все листья в порядке убывания числа появления)

Шаг 2. Пока в очереди больше одного узла, выполняем следующие действия:

Шаг 2а. Удаляем из очереди два узла с наивысшим приоритетом (самыми низкими числами появлений)

Шаг 2б. Создаём новый узел, для которого выбранные на шаге 2а узлы являются наследниками. При этом число появлений нового узла (его вес) полагается равным сумме появлений выбранных на шаге 2а узлов

Шаг 2с. Добавляем узел, созданный на шаге 2б, в очередь с приоритетами.

Единственный узел, который останется в очереди с приоритетами в конце работы алгоритма, будет корневым для построенного двоичного дерева.

Чтобы восстановить кодирующие слова, нам необходимо пройти по рёбрам от корневого узла получившегося дерева до каждого листа. При этом каждому ребру присваивается бит "1", если ребро находится слева, и "0" - если справа (или наоборот).

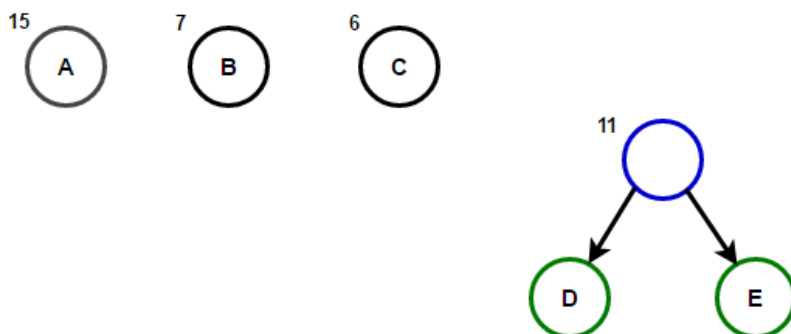
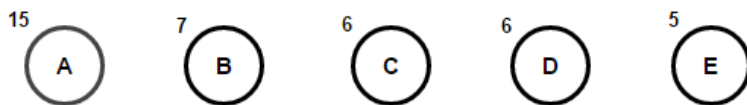
Кодирование Хаффмана

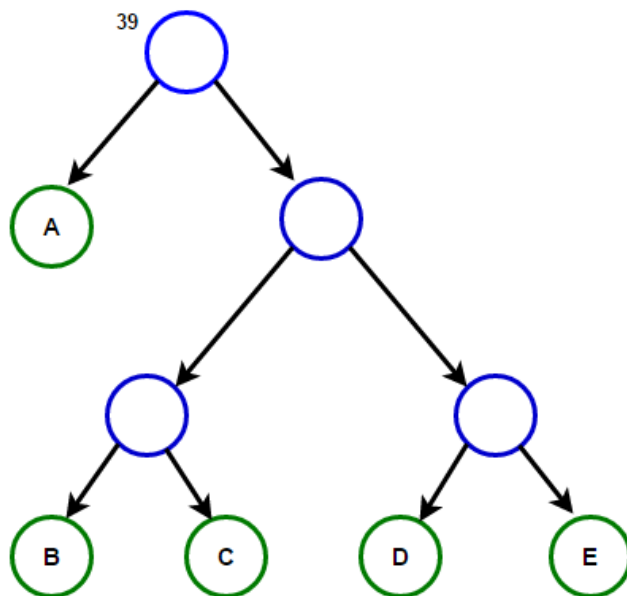
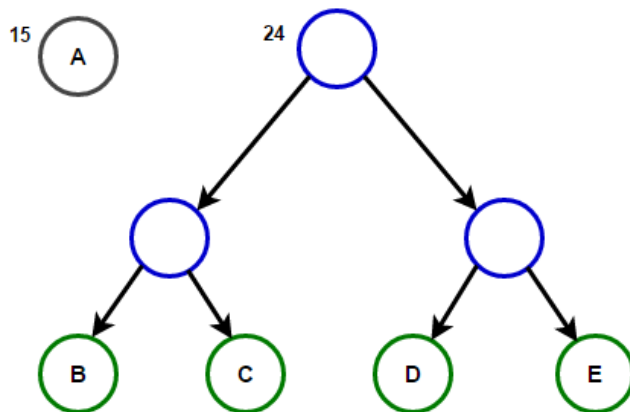
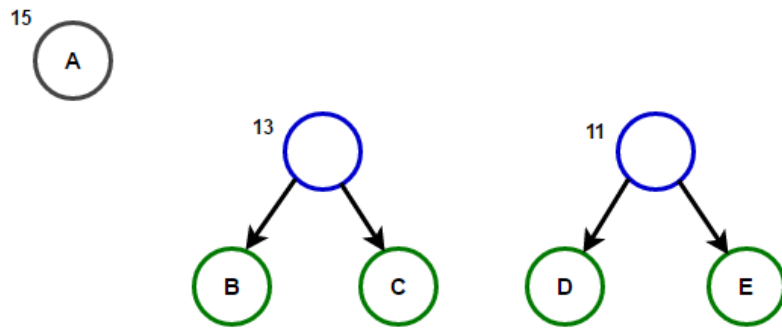
Техника работает, создавая **бинарное дерево** узлов. Узел может быть листовым узлом или внутренним узлом. Изначально все узлы являются листовыми узлами, которые содержат сам персонаж, вес (частоту появления) персонажа. Внутренние узлы содержат вес символов и ссылки на два дочерних узла. По общему соглашению, бит 0 представляет следующий левый дочерний элемент, и немного 1 представляет следующий правый ребенок. Готовое дерево имеет n листовые узлы и $n-1$ внутренние узлы. Рекомендуется, чтобы дерево Хаффмана отбрасывало неиспользуемые символы в тексте, чтобы получить наиболее оптимальную длину кода.

Мы будем использовать **приоритетная очередь** для построения дерева Хаффмана, где узел с наименьшей частотой имеет наивысший приоритет. Ниже приведены полные шаги:

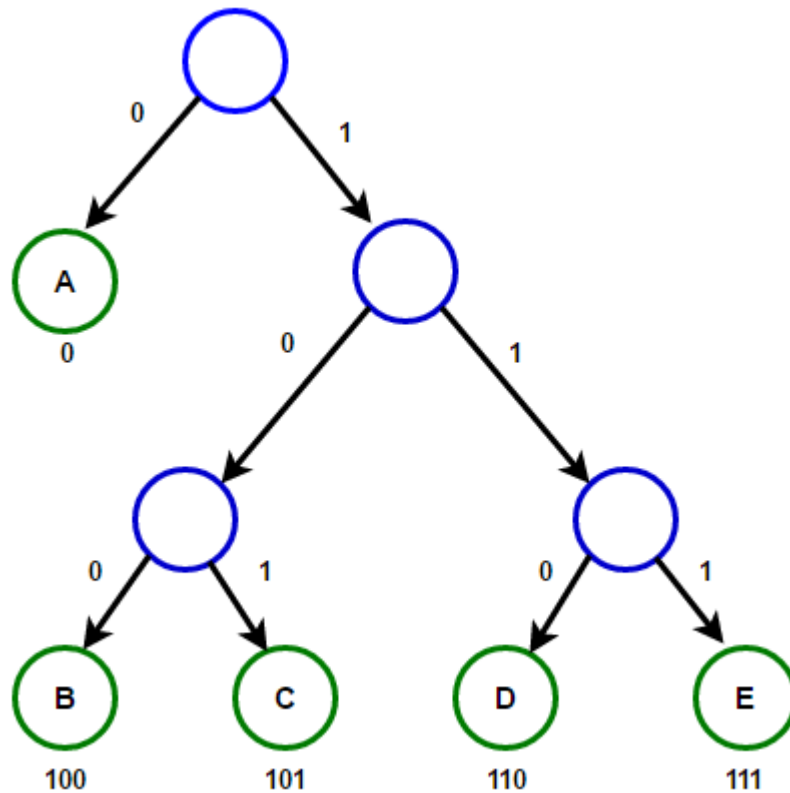
1. Создайте конечный узел для каждого символа и добавьте их в очередь приоритетов.
2. Пока в очереди больше одного узла:
 - Удалите из очереди два узла с наивысшим приоритетом (самой низкой частотой).
 - Создайте новый внутренний узел с этими двумя узлами в качестве дочерних элементов и частотой, равной сумме частот обоих узлов.
 - Добавьте новый узел в очередь приоритетов.
3. Оставшийся узел является корневым узлом, и дерево завершено.

Рассмотрим некоторый текст, состоящий только из 'A', 'B', 'C', 'D', а также 'E' символов, а их частота 15, 7, 6, 6, 5, соответственно. Следующие рисунки иллюстрируют шаги, за которыми следует алгоритм:





Путь от корня к любому конечному узлу хранит оптимальный код префикса (также называемый кодом Хаффмана), соответствующий символу, связанному с этим конечным узлом.



```

1  import java.util.Comparator;
2  import java.util.HashMap;
3  import java.util.Map;
4  import java.util.PriorityQueue;
5
6  // Узел дерева
7  class Node
8  {
9      Character ch;
10     Integer freq;
11     Node left = null, right = null;
12
13     Node(Character ch, Integer freq)
14     {
15         this.ch = ch;
16         this.freq = freq;
17     }
18
19     public Node(Character ch, Integer freq, Node left, Node
20 right)
21     {
22         this.ch = ch;
23         this.freq = freq;
24         this.left = left;
25         this.right = right;
26     }
27 }
28

```

```

29 class Main
30 {
31     // Проходим по дереву Хаффмана и сохраняем коды Хаффмана на
32     карте.
33     public static void encode(Node root, String str,
34                             Map<Character, String> huffmanCode)
35     {
36         if (root == null) {
37             return;
38         }
39
40         // Найден лиственный узел
41         if (isLeaf(root)) {
42             huffmanCode.put(root.ch, str.length() > 0 ? str :
43 "1");
44         }
45
46         encode(root.left, str + '0', huffmanCode);
47         encode(root.right, str + '1', huffmanCode);
48     }
49
50     // Проходим по дереву Хаффмана и декодируем закодированную
51     строку
52     public static int decode(Node root, int index, StringBuilder
53 sb)
54     {
55         if (root == null) {
56             return index;
57         }
58
59         // Найден лиственный узел
60         if (isLeaf(root))
61         {
62             System.out.print(root.ch);
63             return index;
64         }
65
66         index++;
67
68         root = (sb.charAt(index) == '0') ? root.left :
69 root.right;
70         index = decode(root, index, sb);
71         return index;
72     }
73
74     // Вспомогательная функция для проверки, содержит ли дерево
75     Хаффмана только один узел
76     public static boolean isLeaf(Node root) {
77         return root.left == null && root.right == null;
78     }
79
80     // Строит дерево Хаффмана и декодирует заданный входной
81     текст
82     public static void buildHuffmanTree(String text)

```

```

83     {
84         // Базовый случай: пустая строка
85         if (text == null || text.length() == 0) {
86             return;
87         }
88
89         // Подсчитаем частоту появления каждого символа
90         // и сохранить его на карте
91
92         Map<Character, Integer> freq = new HashMap<>();
93         for (char c: text.toCharArray()) {
94             freq.put(c, freq.getOrDefault(c, 0) + 1);
95         }
96
97         // создаем приоритетную очередь для хранения живых узлов
98         // дерева Хаффмана.
99         // Обратите внимание, что элемент с наивысшим
100        // приоритетом имеет наименьшую частоту
101
102         PriorityQueue<Node> pq;
103         pq = new PriorityQueue<>(Comparator.comparingInt(l ->
104        l.freq));
105
106         // создаем конечный узел для каждого символа и добавляем
107        // его
108         // в приоритетную очередь.
109
110         for (var entry: freq.entrySet()) {
111             pq.add(new Node(entry.getKey(), entry.getValue()));
112         }
113
114         // делаем до тех пор, пока в queue не окажется более
115        // одного узла
116         while (pq.size() != 1)
117         {
118             // Удаляем два узла с наивысшим приоритетом
119             // (самая низкая частота) из queue
120
121             Node left = pq.poll();
122             Node right = pq.poll();
123
124             // создаем новый внутренний узел с этими двумя
125            // узлами в качестве дочерних
126             // и с частотой равной сумме обоих узлов
127             // частоты. Добавьте новый узел в очередь
128            // приоритетов.
129
130             int sum = left.freq + right.freq;
131             pq.add(new Node(null, sum, left, right));
132         }
133
134         // `root` хранит указатель на корень дерева Хаффмана
135         Node root = pq.peek();
136

```

```

137         // Проходим по дереву Хаффмана и сохраняем коды Хаффмана
138     на карте
139     Map<Character, String> huffmanCode = new HashMap<>();
140     encode(root, "", huffmanCode);
141
142     // Выводим коды Хаффмана
143     System.out.println("Huffman Codes are: " + huffmanCode);
144     System.out.println("The original string is: " + text);
145
146     // Печатаем закодированную строку
147     StringBuilder sb = new StringBuilder();
148     for (char c: text.toCharArray()) {
149         sb.append(huffmanCode.get(c));
150     }
151
152     System.out.println("The encoded string is: " + sb);
153     System.out.print("The decoded string is: ");
154
155     if (isLeaf(root))
156     {
157         // Особый случай: для ввода типа а, аа, ааа и т. д.
158         while (root.freq-- > 0) {
159             System.out.print(root.ch);
160         }
161     }
162     else {
163         // Снова проходим по дереву Хаффмана и на этот раз
164         // декодируем закодированную строку
165         int index = -1;
166         while (index < sb.length() - 1) {
167             index = decode(root, index, sb);
168         }
169     }
170 }
171
172 // Реализация алгоритма кодирования Хаффмана на Java
173 public static void main(String[] args)
174 {
175     String text = "Huffman coding is a data compression
176 algorithm.";
177     buildHuffmanTree(text);
178 }
179 }

```