

Design:

I started with implementing the Huffman coding. My design started with creating a class for a Huffman node that had a symbol, probability, and a left and right pointer. Next I would build a tree using these nodes, and then use the tree to encode the message. Lastly I would use the tree to decode the message. I started building the tree using a hash map, but quickly found that this was too difficult and fiddly. I ended using a priority queue which allowed me to extract the lowest probability nodes and create a parent node. I continued in this manner until my priority queue only held one node. After creating my tree I then had to encode the message. To do this I walked down the tree till I reached a leaf node appending 1s or 0s depending on which way I went. I then had the code for that specific source symbol which I stored in a hash map. The only problem I ran into here was forgetting to delete a character when I went up the tree. After I had all the symbols encoded it was just a matter of iterating through the message encoding each symbol with their respective codes. Lastly I decoded by following the encoded message down the tree until I reached a leaf, then resetting.

In order to implement the arithmetic coding I used a combination of the algorithm set out in the lecture slides with the one in the suggested reading. I started in a similar manner to the Huffman, by creating a class to hold each symbol with its probability range. I then I had to convert the source symbols with their probabilities into arithmetic symbols. Then encoded the message, and lastly decode it. The most difficult part of the arithmetic coding was attempting to do bitwise calculations. I eventually gave up on this and instead hard coded a high, three quarter, half, and quarter value for the range. With these values I was able to check where the high end and low end of the range fell, for instance if the low end was over or equal to half then a 1 could be added to the encoding. The decoding was a matter of reverse engineering using a buffer of the encoded message. I decided to use the first symbol of the information source as my end of data symbol so I knew where the message ended.

Examples:

with the information source:  
source symbols: a b c d

probabilities: .01 .25 .5 .24

Source Message:

dbbcbdddbdddbdbbbbccbdccbdcdcccddbbbdddbbcbdbbcbdcccbdbbbbcdcdcbdcbbcbcb.  
bdbdbbcbdddbdbbbbcbcdcddbdddbbdbcbbcbdbbbbbcdddbbbdcbbcdcbccbcdbcbbbbcdcbdd  
dbba

Huffman Coding:

Encoded Message:

0000101101000000010000000100001010101110100011010001000111000000010101000000  
0000101101000010110100011110100001010101100000010101000101000101011010101000  
0100001011010000000100001010101101100000001000000000000010100000001101011010  
1010101011000000000010101000101011000000011101100001101010101100010101000000  
0101001

Decoded Message:

dbbcbdddbdddbdbbbbccbdccbdcdcccddbbbdddbbcbdbbcbdcccbdbbbbcdcdcbdcbbcbcb  
bdbdbbcbdddbdbbbbcbcdcddbdddbbdbcbbcbdbbbbbcdddbbbdcbbcdcbccbcdbcbbbbcdcbdd  
dbba

Average Length: 2.0733333333333333

Arithmetic Coding:

Encoded Message:

1100010011001110000010001010111010100101111001110101100100011001000101001011  
110100011001111100111110111101011011110111011111101001100100110110000111011  
0111001001110110000100111100001100111100101001111001111101110010001111110110  
0011001110001011010101100110101000110100110001000000000000000000000000000000

Decoded Message:

dbbcbdddbdddbdbbbbccbdccbdcdcccddbbbdddbbcbdbbcbdcccbdbbbbcdcdcbdcbbcbcb  
bdbdbbcbdddbdbbbbcbcdcddbdddbbdbcbbcbdbbbbbcdddbbbdcbbcdcbccbcdbcbbbbcdcbdd  
dbba

Average Length: 1.98

We can see that the Arithmetic coding outperformed the Huffman coding producing an average length of 1.98 compared to 2.07. With the same information source:

Source Message: cdbdbbcbcbdbbbddccdbbbcbcccbdda

Huffman Coding:

Encoded Message:

10000100001011011010000101010000001100001010110111101000000001

Decoded Message: cdbdbbcbcbdbbbddccdbbbcbcccbdda

Average Length: 2.0

Arithmetic Coding:

Encoded Message:

1010101000001010111100100101110110111110001010001001110001100000000000000000  
00000

Decoded Message: cdbdbbcbcbdbbbddccdbbbcbcccbdda

Average Length: 2.6129032258064515

We can see that as the source message gets smaller the Arithmetic coding is outperformed by the Huffman coding.

Conclusion:

On the whole I found that this practical was difficult because of java and would have possibly been easier in C where more bit control is made available. I feel that if I could do the practical over again I would have read more of the literature on it and drawn out a full game plan before jumping in, as I found myself doing a lot of backtracking.