



| Fixed Length Encoding   |
|---|
| Suppose we have a 4-character alphabet {a, b, c, d}CharacterabcdFixed-Length Codeword00011011 |
| Given the string "abacdaacac", it would be encoded as   |
| a b a c d a a c a c<br>00 01 00 10 11 00 00 10 00 10  |
| The final 20-character binary string would be "00010010110000100010".                         |
| But what if we knew the frequency of the characters in advance?                               |

| Variable Length Encoding   |
|--|
| Using variable length codes  |
| $\begin{array}{ c c c c c c }\hline Character & a & b & c & d \\ \hline Probability & 0.60 & 0.05 & 0.30 & 0.05 \\ \hline Variable-Length Codeword & 0 & 110 & 10 & 111 \\ \hline \end{array}$ |
| Given the string "abacdaacac", it would be encoded as  |
| a b a c d a a c a c<br>0 110 0 10 111 0 0 10 0 10<br>The resulting 17-character string would be "01100101110010010".<br>(savings of 3 bits)  |
| Resulting string is 1.5n compared to 2n, for a savings of 25% in expected encoding length.   |
| $n(0.60 \cdot 1 + 0.05 \cdot 3 + 0.30 \cdot 2 + 0.05 \cdot 3) = n(0.60 + 0.15 + 0.60 + 0.15) = 1.5n.$<br>(Rhodes College   |

| Prefix   | Codes   |  |  |  |  |  |  |
|--|---|--|--|--|--|--|--|
| How to decode variable-length codes?   | CharacterabcdProbability0.600.050.300.05Variable-Length Codeword011010111                               |  |  |  |  |  |  |
| In the variable-length codes given in is a prefix of another (very importa   | n the example above no codeword<br>nt!)   |  |  |  |  |  |  |
| Observe: If two codewords did sha<br>and $b \rightarrow 00101$ , then when we see<br>the first character of the encoded n  | re a common prefix, e.g. a $\rightarrow$ 001<br>00101, how do we know whether<br>nessage is "a" or "b"? |  |  |  |  |  |  |
| Conversely: If no codeword is a prefix of any other, then as soon as we see a codeword appearing as a prefix in the encoded text, then we know that we may decode it |   |  |  |  |  |  |  |
|  | Rhodes5   |  |  |  |  |  |  |





$$B(T) = n \sum_{x \in C} p(x) d_T(x).$$

**Optimal Code Generation**: Given an alphabet C and the probabilities p(x) of occurrence for each character  $x \in C$ , compute a prefix code T that *minimizes* the expected length of the encoded bit-string, B(T).

Rhodes College

n = # of characters in the encoded string





| Huffman Code Construc    |      |      |  |
|--------------------------|------|------|--|
|                          | Char | Freq |  |
| Character count in text. | E    | 125  |  |
|                          | Т    | 93   |  |
|                          | A    | 80   |  |
|                          | 0    | 76   |  |
|                          | I    | 73   |  |
|                          | N    | 71   |  |
|                          | 5    | 65   |  |
|                          | R    | 61   |  |
|                          | Н    | 55   |  |
|                          | L    | 41   |  |
|                          | D    | 40   |  |
|                          | С    | 31   |  |
|                          | U    | 27   |  |





























25



Recall that the cost of any encoding tree T is  $B(T) = n \sum_{x \in T} p(x) d_T(x)$ .

Need to show that any tree that differs from the one constructed by Huffman's algorithm can be converted into one that is equal to Huffman's tree without increasing its cost

The key is showing that the greedy choice is always the proper one to make (or at least it is as good as any other choice).

Our approach is based a few observations.

- 1. The Huffman tree is a full binary tree, meaning that every internal node has exactly two children.
- 2. The two characters with the lowest probabilities will be siblings at the maximum depth in the tree.







