

CS 512: The Cost of Information

16:198:512

Instructor: Wes Cowan

Consider the following text:

Arma virumque cano, Troiae qui primus ab oris
 Italiam, fato profugus, Laviniaque venit
 litora, multum ille et terris iactatus et alto
 vi superum saevae memorem Iunonis ob iram;
 multa quoque et bello passus, dum conderet urbem,
 inferretque deos Latio, genus unde Latinum,
 Albanique patres, atque altae moenia Romae.

What is the information content of this paragraph? Unless you speak Latin it is tempting to say not much. However, we can address the question in almost a context-free fashion by asking the following question: how much data would I have to send to you, in order to communicate this passage to you? In standard ASCII encoding, each character is represented by 8 bits, which means (ignoring new lines) we could express this paragraph in 314 symbols * 8 bits per symbol = 2512 bits. But can we do better? Noticing that there are only 28 characters in the whole text, we could make do with a 5 bit encoding (there are a lot of symbols we wouldn't have to worry about) which would potentially reduces us to a total 1570 bits.

Scanning through the text (again, ignoring new lines), there are 28 characters occurring with the following frequencies:

	47	e	31	a	27	u	25	i	21	t	20	m	17
o	17	r	17	s	14	n	12	l	11	,	9	q	8
p	5	v	5	b	5	d	4	f	3	c	3	L	3
A	2	I	2	g	2	.	1	;	1	R	1	T	1

From this table, we see that the space character occurs vastly more often than the letter *T*. If we could come up with an encoding that used fewer bits for the space character than it did for the letter *T*, (and generally required fewer bits for more frequent symbols) we could potentially reduce the total number of bits required to communicate the text from one person to another. Consider the following bitwise encoding:

	01	e	011	a	0111	u	01111	i	011111	t	0111111	m	01111111
o	001	r	0011	s	00111	n	001111	l	0011111	,	00111111	q	001111111
p	0001	v	00011	b	000111	d	0001111	f	00011111	c	000111111	L	0001111111
A	00001	I	000011	g	0000111	.	00001111	;	000011111	R	0000111111	T	00001111111

One of the nice things about this encoding is that it is 'prefix free', that is, no encoding occurs as the prefix of another encoding. What this means is that when the text is presented in encoded form, the bits can naturally be partitioned and the symbols recovered. Encoding the original text (ignoring newlines) gives:

```
00001001101111111011101000110111110011011110111111100111111101111
0110100011111101110011110010011111101000011111110011001011111011
101101001111111011110111110100010011011111011111101111100111010111
00011101001001101111100111010000110111111011100111110111101110111
111100111111010001111101110111111001010001001100100011111011110000
11101111001110011111101000111111011100011011110011110111110111001
```

```

1111110111101101000110110011110111110111111010011111011111011111100
100110111001111110101111111011110011111011111101111110111111101011111
001111100111110110101101111110101111110110011001101111100111010111
110111000111111011111011101111110111100111010110111111010111001111
10111111001010001101111101001110111100010110011011110111111010011
1011101100011011101101011111101101111110010011011011111110100001
101111001111001001111011111001110100100011101011111001101110111111
100001111101011111101111001111101111110111010011111110111100100111
111101111011010110111111010001110110011111001111100101000101110011
10011101111001110011111101000111011110111111010001111110010011110
001111011001101101111110101111001100011101101111110011111101011111
001111000111110110011001101101111110011111101111011010001111011001
001110100011111101110111111011111001001111110100001110110011110111
1001110101111001111000111101101000111111101111101111110011111011
1101111111001111110100001001111100011101110011110111110011111110111
101101000101110111110011011001110011111101011101111110011111110111
1011010111001111101111110111011010111111100101100111101111101110100
001111110010111111011101100001111
    
```

Looking at the table, it is easy to parse through and see that the first character is an *A*, and see where the second character starts. The above message is a total of 1562 bits, which is an improvement (though not by much) on the naive 5-bit encoding we could use. Can we do better? Enter Huffman encoding.

Huffman encoding proceeds the following way (see the text for a more graphical / algorithmic description): first, take the two symbols with the lowest frequencies, in this case *T* and *R*. We build a prefix-free encoding starting by declaring that the *last digit of the encoding of T will be 0, and the last digit of the encoding of R will be 1*. But how do determine the *rest* of the encoding? Imagine merging *T* and *R* into a meta-symbol *TR*, which occurs with the combined frequency $1 + 1 = 2$. Replacing *T* and *R* with this meta-symbol *TR*, we can ask “*what should the encoding of TR be?*” This is a smaller version of the same problem we are originally presented with:

	47	e	31	a	27	u	25	i	21	t	20	m	17
o	17	r	17	s	14	n	12	l	11	,	9	q	8
p	5	v	5	b	5	d	4	f	3	c	3	L	3
A	2	I	2	g	2	TR	2	.	1	;	1		

Looking at this table, we can take the two least frequent symbols, *.* and *,*, and declare that the last digit of the encoding of *,* should be 0 and the last digit of the encoding for *.* should be 1. Again, we merge these two into a meta-symbol *;*, and reinsert it into the table with its combined frequency $1 + 1 = 2$.

	47	e	31	a	27	u	25	i	21	t	20	m	17
o	17	r	17	s	14	n	12	l	11	,	9	q	8
p	5	v	5	b	5	d	4	f	3	c	3	L	3
A	2	I	2	g	2	TR	2	;	2				

At this point, we look at the two least frequent symbols, the meta-symbol *;*, and the meta-symbol *TR* and declare the last digit of their encodings to be 0 and 1 respectively. Note that at this point, we have that the last *two* digits of the encoding for *T* are 10, the last two digits of *R* are 11, the last two digits of *,* are 00, and the last two digits of

. are 01. We again merge, constructing a meta-symbol ;*TR* with frequency $2 + 2 = 4$. This process iterates until all symbols and meta symbols have been merged into the final two, which (by my iteration) yields **sgI;TRqeromdAL**, which will be encoded with a 0, and **tbvpcfiln ua** which will be encoded with a 1. Unpacking all this yields the full encoding for each character as

	110	e	001	a	1111	u	1110	i	1010	t	1000	m	0110
o	0101	r	0100	s	0000	n	10111	l	10110	,	01111	q	00011
p	100110	v	100101	b	100100	d	011100	f	1001111	c	1001110	L	0111011
A	0111010	I	0001001	g	0001000	.	00010101	;	00010100	R	00010111	T	00010110

This leads to an encoded message of:

```
0111010010001101111110100101101001001110011000011111
00011101001110111101110101011111000010110010001011
0101111001110000111110101011010011001001010011011100
000110111110010011001010100101000001100001001100011
1110110101011110110011111101001111111100001011101001
10010001011001111111000010001110000001111110011101111
11100101101010111101011110001111100011101001010011011
11010100011010110101010000101010011110111111001101110
10110100011100110110101010110101100011100011000110100
00010100010010100000110101011111001110100011111000111
0000011000110001101111011010000101110100101101011000
00111010011000101001110011011000001111001100101111100
11100110001011001010100001011011000010011110101110101
10111101000001100101100100110101001001111011000010100
110011011101011010001111110000111110010100011111000111
00011000110100100001101101011001011101001101111000000
001110000001111110011100111001101101001110010110111011
10000101000011000110111001001001000010110011111101010
101111001111001010001000011000000111110001110011100001
0101000011001110111111000101001010111111000010000011
011111100000110111010111011100001110011101111111000101
010111111001100111111001110101011010010011111011110100
001111100011101001101111100001000010000011111101111100
000011111000111011111011010001111001110011001010011011
1101011111100001011101010110111100100010101
```

This may not look like much by itself, but it is worth noting that this comes out to 1317 bits, a 52% reduction in size over the original ASCII encoding. That being so, it is worth asking - what is the best we could hope to achieve? Is there a lower bound on how small we could go? If we were paying attention to context, we might compress the original text as “Virgil’s Aeneid, Book 1, Card 1.” This is definitely an impressive compression, but unfortunately this method doesn’t generalize to compressing anything that is not in Virgil’s Aeneid, or anything that is not already written somewhere else in an easily referenced form. In comparison, the encoding constructed above could be used to compress and transmit *any* message (subject to it containing only the same 28 symbols as our original text here). Taking the view of messages as sequences of independent symbols occurring with some frequency or probability (divorced from a larger cultural understanding), Shannon’s Information Theory tells us that there is a lower bound

on the amount of information a message contains, summarized in terms of the *entropy*, or effective average bits per symbol. If a corpus contains N symbols, with symbol i occurring with frequency p_i , the bit entropy is defined to be

$$H(\text{corpus}) = - \sum_{i=1}^N p_i \log_2 p_i. \quad (1)$$

For our ‘corpus’ as above, we get $H \approx 4.1664$ bits per symbol, for a total over 314 symbols of 1308.25 bits. As such, the Huffman encoding is essentially optimal, the discrepancy between the observed 1317 bits and the theoretical bound of 1308.25 bits occurring being largely due to the frequencies not working out to exact powers of two. Note additionally, this demonstrates that the naive 5-bit / symbol encoding improvement on ASCII encoding is pretty close to optimal.

Questions:

- 1) We could also consider compressing the text by ‘chunking’ or blocking it into pairs of characters, [Ar][ma][v]... , and then treating each pair as a symbol, and encoding over this expanded set of symbols. Do that in this case - is the compression worse, better, or effectively no different, in terms of total number of bits required for the message? What is the bitwise entropy for the message when symbols are taken as pairs? In general, what would lead to blocking by pairs giving a better compression than individual symbols, and what would lead to blocking by pairs giving worse compression?
- 2) Suppose that you have a corpus over an N -symbol alphabet, each symbol occurring with frequency or probability p_i . How many possible pairs of symbols are there, and for any symbol pair, what is its frequency in this corpus? *Assume symbols occur independently with the indicated frequencies.*
- 3) In Question (2), how does the bitwise entropy of the individual symbol frequencies compare to the bitwise entropy of the pairwise symbol frequencies? What can you conclude from this?
- 4) Argue that for a given set of symbols, there is no encoding that can universally compress all possible messages in those symbols. Why is this not a problem, i.e., why do we still pursue these encodings when they cannot be universally effective?