

Math 5251 Huffman Coding (§3.4)

It turns out that given the source word probabilities (p_1, \dots, p_m) for $W = \{w_1, \dots, w_m\}$, we can easily find an n -ary encoding $f: W \rightarrow \Sigma^*$ that achieves the minimum for $\text{avg length}(f)$, via **Huffman coding**.

Let's

- describe the binary case first,
($\Sigma = \{0, 1\}$)
- prove that it achieves the minimum,
- then explain how to modify it for n -ary.

Binary Huffman encoding algorithm:

Assume by re-indexing that

$$p_1 \geq p_2 \geq \dots \geq p_{m-2} \geq p_{m-1} \geq p_m$$

and **recursively** define $f: W \rightarrow \{0,1\}^*$ by induction on m :

If $m=2$, (BASE CASE) so $W = \{w_1, w_2\}$ encode $f(w_1) = 0$
probabilities p_1, p_2 $f(w_2) = 1$

If $m > 2$, build a Huffman encoding for a source $W' = \{w'_1, w'_2, \dots, w'_{m-2}, w'_{m-1}\}$ with probabilities $\{p_1, p_2, \dots, p_{m-2}, p_{m-1} + p_m\}$ and then tack on an extra 0 to $f(w'_{m-1})$ and an extra 1

$$\text{i.e. } f(w_i) = \begin{cases} f(w'_i) & \text{if } i = 1, 2, \dots, m-2 \\ f(w'_{m-1})0 & \text{if } i = m-1 \\ f(w'_{m-1})1 & \text{if } i = m \end{cases}$$

Usually this is visualized via binary Huffman trees, reading code words as paths from **root** to **leaves**...

EXAMPLES

(1) $W = \{A, B, C, D\}$

probabilities $\frac{1}{2} \geq \frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{10}$
 $p_1 \quad p_2 \quad p_3 \quad p_4$
 add these, giving

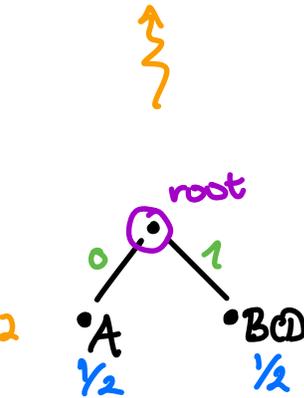
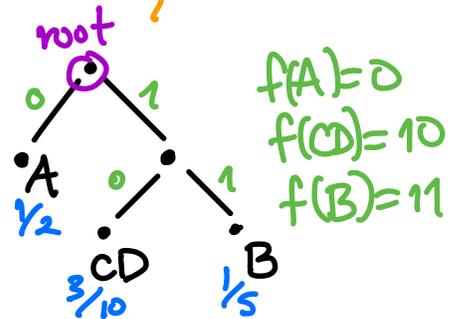
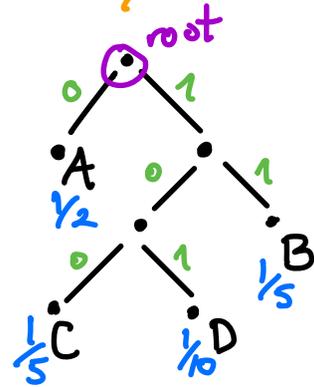
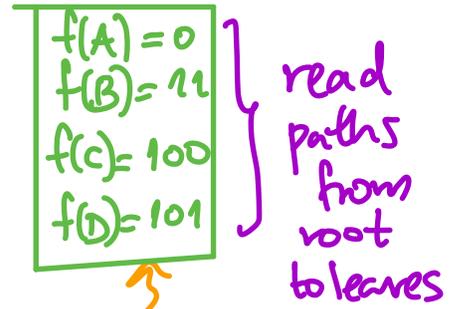
$W' = \{A, CD, B\}$

$\frac{1}{2} \geq \frac{3}{10} \geq \frac{1}{5}$
 add these, giving

$W'' = \{A, BCD\}$

$\frac{1}{2} \geq \frac{1}{2}$

base case $n=2$



$f(A) = 0$
 $f(BCD) = 1$

(2) If some p_i coincide (or their sums coincide), the Huffman encoding may not be unique, e.g.

$$W = \{A, B, C, D, E\}$$

$$\frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5}$$

$$W' = \{DE, A, B, C\}$$

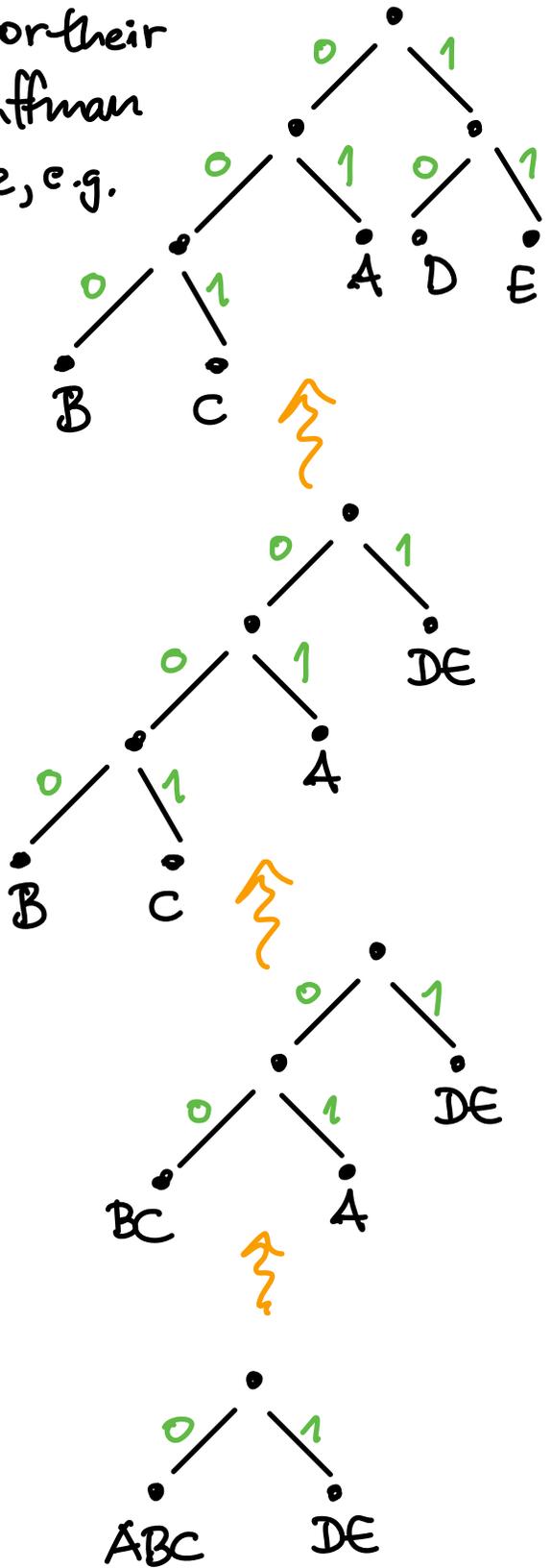
$$\frac{2}{5} \geq \frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5}$$

$$W'' = \{DE, BC, A\}$$

$$\frac{2}{5} \geq \frac{2}{5} \geq \frac{1}{5}$$

$$W''' = \{ABC, DE\}$$

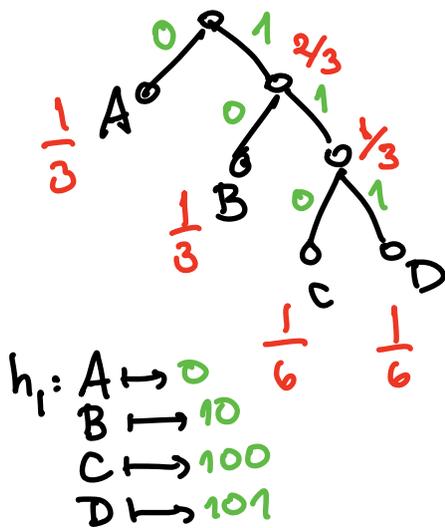
$$\frac{3}{5} \geq \frac{2}{5}$$



BETTER EXAMPLE of non-uniqueness.

$W = \{A, B, C, D\}$ has two possible binary Huffman tree structures, having different codeword lengths (but necessarily same avg length):

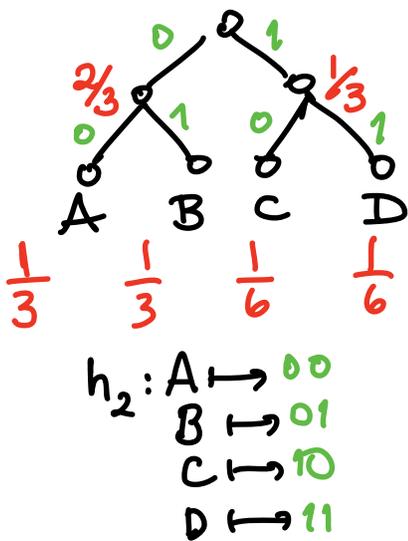
probos $\frac{1}{3} \frac{1}{3} \frac{1}{6} \frac{1}{6}$



$$(l_1, l_2, l_3, l_4) = (1, 2, 3, 3)$$

$$\text{avg length}(h_1) =$$

$$\frac{1}{3} \cdot 1 + \frac{1}{3} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{1}{6} \cdot 3$$
$$= \frac{2+4+3+3}{6} = 2$$



$$(l_1, l_2, l_3, l_4) = (2, 2, 2, 2)$$

$$\text{avg length}(h_2) =$$

$$\frac{1}{3} \cdot 2 + \frac{1}{3} \cdot 2 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 2$$
$$= 2$$

THEOREM Let $W = \{\omega_1, \dots, \omega_m\}$ have probabilities $\{P_1, \dots, P_m\}$ and $h: W \rightarrow \{0,1\}^*$ any Huffman encoding.

Then (a) h is prefix, so u.d., and

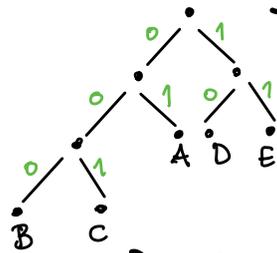
(b) for any u.d. encoding $f: W \rightarrow \{0,1\}^*$

$$\text{avg length}(h) \leq \text{avg length}(f)$$

(so h achieves the minimum bounded in Shannon's Thm.)

EXAMPLE This Huffman encoding has

$W = \{A, B, C, D, E\}$
 $\frac{1}{5}, \frac{2}{5}, \frac{1}{5}, \frac{2}{5}, \frac{1}{5}$



$A \mapsto 01$
 $B \mapsto 000$
 $C \mapsto 001$
 $D \mapsto 10$
 $E \mapsto 11$

with lengths $(l_1, l_2, l_3, l_4, l_5) = (2, 2, 2, 3, 3)$.

Why can't we find something **shorter**, like $(2, 2, 2, 2, 3)$?

proof of THEOREM:

For (a), note that each Huffman codeword $f(\omega)$ is the labels on a path from root to a leaf in the tree. So $f(\omega)$ can't be a prefix of another $f(\omega')$, else the path from the root **continues lower**, so it wasn't stopping at a **leaf** to read $f(\omega)$.

For (b), assume that $f: W \rightarrow \{0,1\}^*$ is a u.d. encoding achieving the **minimum** of $\text{avglength}(f)$ among **all u.d. encodings**. We'll show $\text{avglength}(h) \leq \text{avglength}(f)$ in several steps.

STEP 1: We can **assume f is prefix**, not just u.d., because of the Kraft-McMillan Theorems: the lengths (l_1, \dots, l_m) for $f(w_1), \dots, f(w_m)$ satisfy $\sum_{i=1}^m \frac{1}{n^{l_i}} \leq 1$ and hence \exists a prefix code with the same lengths.

STEP 2: We can assume after re-indexing that if $p_1 \geq p_2 \geq \dots \geq p_{m-2} \geq p_{m-1} \geq p_m$ then f has **$l_1 \leq l_2 \leq \dots \leq l_{m-2} \leq l_{m-1} \leq l_m$** . Otherwise, if $l_i > l_{i+1}$, swap images $f(w_i), f(w_{i+1})$ of w_i, w_{i+1} creating a new u.d. f with smaller $\text{avglength}(f) = \sum_{i=1}^m p_i l_i$.

STEP 3: We can **assume $l_{m-1} = l_m$** , otherwise if $l_{m-1} < l_m$ then we can drop the last letter of $f(w_m)$ without ruining the prefix property (**Why?**), and making $\text{avglength}(f)$ smaller.

STEP 4: We can assume \exists some $i \leq m-1$ such that $f(w_i)$ and $f(w_m)$ have same length $l_i = l_m$ and differ only in their last digit:

$$f(w_i) = a_1 a_2 \dots a_{l-1} 0$$

$$f(w_m) = a_1 a_2 \dots a_{l-1} 1$$

(In which case, re-index so that $i = m-1$).

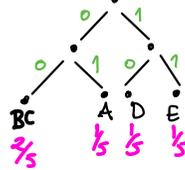
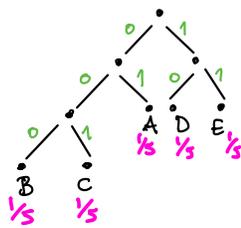
This is because otherwise, we could again drop the last letter of $f(w_m)$ without ruining the prefix property (why?), but reducing $\text{avg length}(f)$.

LAST (INDUCTIVE) STEP:
on m

Create the **smaller Huffman code** $h': W' \rightarrow \{0,1\}^*$ for the source with probabilities $p_1, p_2, \dots, p_{m-2}, p_{m-1} + p_m$ by removing the final 0 from $h(w_{m-1})$ and 1 from $h(w_m)$

$$W = \{A, B, C, D, E\}$$

$$\frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5} \geq \frac{1}{5}$$



Similarly create the **smaller prefix code** $f': W' \rightarrow \{0,1\}^*$ for that same source W' by removing the final 0 from $f(w_{m-1})$ and 1 from $f(w_m)$.

Note how avg length for h and h' relate:

if the Huffman codewords have lengths $\hat{l}_1 \geq \dots \geq \hat{l}_{m-2} \geq \hat{l}_{m-1} = \hat{l}_m$,

$$\text{avglength}(h) = p_1 \hat{l}_1 + \dots + p_{m-2} \hat{l}_{m-2} + \underbrace{p_{m-1} \hat{l}_{m-1} + p_m \hat{l}_m}_{= (p_{m-1} + p_m) \hat{l}_m}$$

$$\text{avglength}(h') = p_1 \hat{l}_1 + \dots + p_{m-2} \hat{l}_{m-2} + (p_{m-1} + p_m) (\hat{l}_m - 1)$$

$$\Rightarrow \boxed{\text{avglength}(h) = \text{avglength}(h') + p_{m-1} + p_m}$$

Similarly,

$$\boxed{\text{avglength}(f) = \text{avglength}(f') + p_{m-1} + p_m}$$

This lets us prove $\text{avglength}(h) \leq \text{avglength}(f)$

by induction on $m = |W|$, since it's easy to check in the base case where $m=2$ (so $h(A)=0$, $h(B)=1$)

and then in the inductive step, use

$$\text{avglength}(h') \leq \text{avglength}(f')$$

together with the two boxed facts above. \square

It's easy to modify Huffman coding for an n -ary alphabet $\Sigma = \{0, 1, 2, \dots, n-1\}$:

the Huffman trees are n -ary and built by

grouping $p_1 \geq p_2 \geq \dots \geq p_{n-m} \geq \underbrace{p_{n-m+1} \geq \dots \geq p_{m-1} \geq p_m}$

$$p_1 \geq p_2 \geq \dots \geq p_{n-m} \geq \sum_{i=n-m+1}^m p_i \quad \text{in } W'$$

The only issue is that n -ary trees have their number of leaves $\equiv 1 \pmod{n-1}$

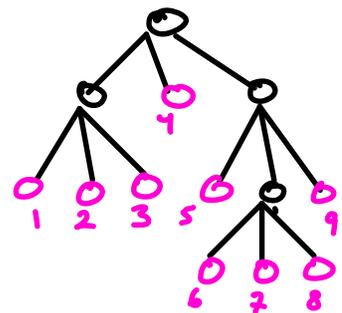
i.e. remainder of 1 on division by $n-1$.

So one pads $p_1 \geq \dots \geq p_m \rightsquigarrow p_1 \geq \dots \geq p_m \geq 0 \geq \dots \geq 0$
with zeroes to make $M \equiv 1 \pmod{n-1}$. $\underset{=}{\parallel}$
 p_m

EXAMPLE

$n=3$ Ternary trees have

number of leaves $\equiv 1 \pmod{2}$
i.e. **odd**

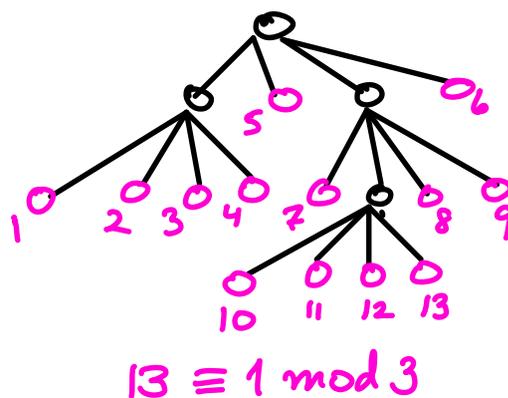


$9 \equiv 1 \pmod{2}$
odd

EXAMPLE

$n=4$

4-ary trees have
number of leaves $\equiv 1 \pmod{3}$



EXAMPLE Morse code is a ternary and prefix
code $f: W = \{A, B, C, \dots, z\} \rightarrow \{0, -, \text{space}\}^* = \Sigma^*$
 $m=26$ $n=3$

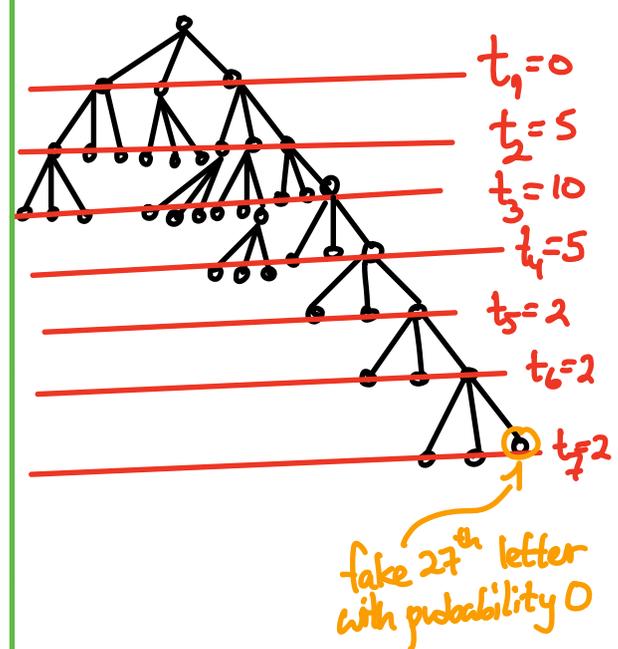
How well does a ternary Huffman code $h: W \rightarrow \{0, 1, 2\}$
beat its avg length?

Since $n=26 \not\equiv 1 \pmod{3}$, need to add
an extra fake 27th letter with probability $p_{27}=0$,
then use a computer to build a ternary Huffman tree...

Letter English Probability Morse code lengths (with space) Ternary Huffman code lengths

:	E	0.12702	2	2
:	T	0.09056	2	2
:	A	0.08167	3	2
:	O	0.07507	3	2
:	I	0.06966	3	2
:	N	0.06749	3	3
:	S	0.06327	4	3
:	H	0.06094	4	3
:	R	0.05987	4	3
:	D	0.04253	4	3
:	L	0.04025	4	3
:	C	0.02782	4	3
:	U	0.02758	4	3
:	M	0.02406	4	3
:	W	0.0236	5	3
:	F	0.02228	5	4
:	G	0.02015	5	4
:	Y	0.01974	5	4
:	P	0.01929	5	4
:	B	0.01492	5	4
:	V	0.00978	5	5
:	K	0.00772	5	5
:	J	0.00153	5	6
:	X	0.0015	5	6
:	Q	0.00095	5	7
:	Z	0.00074	5	7

Ternary Huffman code tree structure:



avg length (h) = 2.7

Morse code (with final space) has length tallies

$(t_1, t_2, t_3, t_4, t_5, t_6, t_7)$
 $= (0, 2, 4, 8, 12, 0, 0)$

avg length₂ (f) = 3.41

REMARK

Although a Huffman encoding achieves the minimum for avg length (f) among u.d. codes, it may not get as low as Shannon's $\frac{H(W)}{\log_2(n)}$

lower bound. But one way to improve it is by grouping source words $W = \{w_1, \dots, w_n\}$ into sequences $W^{(l)} = \{(w_{i_1}, w_{i_2}, \dots, w_{i_l}) : w_i \in W\}$ sent l at a time, called the l^{th} extension of W , with $P(w_{i_1}, w_{i_2}, \dots, w_{i_l}) = p_{i_1} \cdot p_{i_2} \cdots p_{i_l}$

EXAMPLE $W = \{A, B\}$

$$\text{has } H(W) = \frac{3}{4} \log_2\left(\frac{4}{3}\right) + \frac{1}{4} \log_2(4) \approx 0.811278$$

and binary Huffman encoding $f(A) = 0$
 $f(B) = 1$

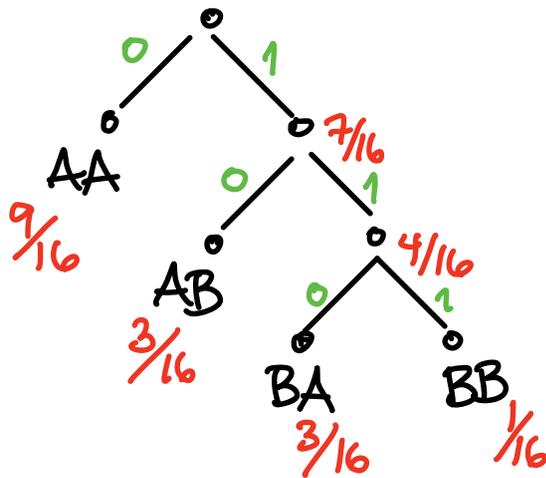
$$\text{with avg length}(f) = \frac{3}{4} \cdot 1 + \frac{1}{4} \cdot 1 = 1 \quad \left(> \begin{matrix} 0.811278 \\ = H(W) \end{matrix} \right)$$

But its 2nd extension

$$W^{(2)} = \{ AA, AB, BA, BB \}$$

$$\begin{array}{cccc} \frac{3}{4} \cdot \frac{3}{4} & \frac{3}{4} \cdot \frac{1}{4} & \frac{1}{4} \cdot \frac{3}{4} & \frac{1}{4} \cdot \frac{1}{4} \\ = \frac{9}{16} & = \frac{3}{16} & = \frac{3}{16} & = \frac{1}{16} \end{array}$$

has binary Huffman encoding as shown:



$$\text{so } \text{avglength}(f) = \frac{9}{16} \cdot 1 + \frac{3}{16} \cdot 2 + \frac{3}{16} \cdot 3 + \frac{1}{16} \cdot 3$$

$$\quad \quad \quad \underbrace{\quad}_{l(0)} \quad \quad \underbrace{\quad}_{l(10)} \quad \quad \underbrace{\quad}_{l(110)} \quad \quad \underbrace{\quad}_{l(111)}$$

$$= \frac{27}{16} = 1.6875$$

But it makes sense to divide this by 2, since we're sending 2 words at a time:

$$\frac{\text{avglength}(f)}{2} = \frac{27}{32} = 0.84375, \text{ much closer to } H(W) \approx 0.81278$$

In fact, its 3rd extension

$$W^{(3)} = \{AAA, AAB, ABA, BAA, ABB, BAB, BBA, BBB\}$$

$$\text{probs } \frac{27}{64} \quad \frac{9}{64} \quad \frac{9}{64} \quad \frac{9}{64} \quad \frac{3}{64} \quad \frac{3}{64} \quad \frac{3}{64} \quad \frac{1}{64}$$

gets amazingly close: $\frac{\text{avglength}(f)}{3} = 0.811278$
matching to 6 digits!

It's not hard to show this version of
Shannon's Noiseless Coding Thm:

(Roman Thm 2.3.4)

THEOREM: The l^{th} extension $W^{(l)}$ of a source W
has entropy $\frac{H(W^{(l)})}{l} = H(W)$,

and among all n -ary c.d. encodings $f: W^{(l)} \rightarrow \Sigma^*$,
the ones achieving minimum $\text{avglength}(f)$ have

$$\frac{H(W)}{\log_2(n)} \leq \frac{\text{avglength}(f)}{l} \leq \frac{1}{l} + \frac{H(W)}{\log_2(n)}$$

can be made
smaller by picking
 l larger