Information Theory

05 Symbol Codes



TI 2020/2021

Notice

Author

João Moura Pires (jmp@fct.unl.pt)

This material can be freely used for personal or academic purposes without any previous authorization from the author, provided that this notice is maintained/kept.

For commercial purposes the use of any part of this material requires the previous authorization from the author.



Bibliography

Many examples are extracted and adapted from:



Information Theory, Inference, and Learning Algorithms

Cambridge University Press, 2003

Information Theory, Inference, and Learning Algorithms David J.C. MacKay 2005, Version 7.2

- And some slides were based on lain Murray course
 - http://www.inf.ed.ac.uk/teaching/courses/it/2014/



Table of Contents

- Symbol codes
- What limit is imposed by unique decodeability?
- What's the most compression that we can hope for?
- How much can we compress?
- Huffman coding
- Disadvantages of the Huffman code
- Summary



Information Theory

Symbol codes



Symbol codes

- Variable-length symbol codes,
 - Encode **one source symbol at a time** (instead of encoding huge strings of *N* source symbols);
 - These codes are **lossless**: they are guaranteed to compress and decompress without any errors;
 - There is a chance that the codes may sometimes produce encoded strings longer than the original source string.

The idea is that we can achieve compression, **on average**, by **assigning shorter encodings** to the **more probable outcomes** and longer encodings to the less probable



Key issues

What are the implications if a symbol code is lossless?

If some codewords are shortened, by how much do other codewords have to be lengthened?

Making compression practical.

How can we ensure that a symbol code is easy to decode?

Optimal symbol codes.

- How should we assign codelengths to achieve the best compression,
- What is the best achievable compression?

Source coding theorem (symbol codes)



A (binary) symbol code

Notation on Alpahabets

- A^N . denotes the set of ordered *N*-tuples of elements from the set *A*, i.e., **all strings of length** *N*.
- A+ will denote the set of **all strings of finite length** composed of elements from the set A.
- Examples:
 - $[0,1]^3 = \{000,001,010,011,100,101,110,111\}.$
 - $\{0,1\}^+ = \{0,1,00,01,10,11,000,001,\ldots\}.$



A (binary) symbol code

- A (binary) symbol code *C* for an ensemble *X* is a mapping from the range of *x*, $A_X = \{a_1, ..., a_I\}$, to $\{0, 1\}^+$.
 - c(x) will denote the codeword corresponding to x
 - *l*(*x*) will denote its length, with $l_i = l(a_i)$.
- The *extended code* C^+ is a **mapping** from A_X^+ to $\{0,1\}^+$ obtained by concatenation, without punctuation, of the corresponding codewords: $c^+(x_1x_2...x_N) = c(x_1)c(x_2)...c(x_N)$
- Basic requirements for a useful symbol code
 - Any encoded string must have a **unique decoding**
 - The symbol code must be **easy to decode**
 - The code should achieve **as much compression as possible**.



A (binary) symbol code - example

- A (binary) symbol code C_0 for an ensemble X $A_X = \{a, b, c, d\}$
 - $P_X = \{1/2, 1/4, 1/8, 1/8\}$
 - c(x) will denote the codeword corresponding to x
 - *l*(*x*) will denote its length, with $l_i = l(a_i)$.

Using the extended code, we may encode acdbac as

 $c^{+}(acdbac) = 100000100001010010000010$ a c d b a c

	a_i	$c(a_i)$	l_i
	a	1000	4
C_0 :	b	0100	4
	С	0010	4
	d	0001	4

Any encoded string must have a unique decoding

• A code C(X) is **uniquely decodeable**

If, under the extended code C^+ , no two distinct strings have the same encoding, i.e.,

 $\forall x, y \in A_X^+, x \neq y \Longrightarrow c^+(x) \neq c^+(y)$

The example is a uniquely decodeable code $c^+(acdbac) = 100000100001010010000010$ $a^+c^+d^-b^+a^+c^-$

How to prove for this example?

	a_i	$c(a_i)$	l_i
	a	1000	4
C_0 :	b	0100	4
	С	0010	4
	d	0001	4



The symbol code must be easy to decode

- A symbol code is easiest to decode if it **is possible to identify the end of a codeword as soon as it arrives**,
 - which means that **no codeword can be a** *prefix* **of another codeword**

- A word c is a prefix of another word d
 - if there exists a tail string *t* such that the concatenation *ct* is identical to *d*.
- Example
 - 1 is prefix of 100



10 is prefix of 100



The symbol code must be easy to decode

A symbol code is called a prefix code

if no codeword is a prefix of any other codeword.

- A prefix code is also known as an *instantaneous* or *self-punctuating* code, because an encoded string can be decoded from left to right without looking ahead to subsequent codewords.
- The end of a codeword is immediately recognizable.
- A prefix code is uniquely decodeable.



Prefix codes - examples

The code $C_1 = \{0, 101\}$ is a prefix code because 0 is not a prefix of 101, nor is 101 a prefix

of 0.



 C_1 is an *incomplete* code

The code $C_2 = \{1, 101\}$ is not a prefix code because 1 is a prefix of 101.





The code should achieve as much compression as possible

The expected length L(C, X) of a symbol code C for an ensemble X is

$$L(C, X) = \sum_{x \in A_X} P(x)l(x)$$
$$L(C, X) = \sum_{i=1}^{I} p_i l_i \quad \text{where} \quad I = |A_X|$$

$$l_i = h_i = \log_2(1/p_i)$$

 $A_X = \{a, b, c, d\}$ H(X) = 1.75 bits

 $P_X = \{1/2, 1/4, 1/8, 1/8\}$





	a_i	$c(a_i)$	p_i	$h(p_i)$	l_i
	a	0	1/2	1.0	1
	b	10	$1/_{4}$	2.0	2
	С	110	1/8	3.0	3
L(C,X) = 1.75 bits	d	111	1/8	3.0	3

The code should achieve as much compression as possible

A (binary) symbol code C₆ for an ensemble X

 $A_X = \{a, b, c, d\}; P_X = \{1/2, 1/4, 1/8, 1/8\}$ H(X) = 1.75 bits

		C_6 .		
a_i	$c(a_i)$	p_i	$h(p_i)$	l_i
а	0	$1/_{2}$	1.0	1
b	01	$1/_{4}$	2.0	2
С	011	1/8	3.0	3
d	111	1/8	3.0	3

$$C_6$$
:

$$L(C,X) = 1.75 bits$$

Is C_6 a prefix code? No

Is C_6 a uniquely decodeable? Yes

Certainly isn't easy to decode



The code should achieve as much compression as possible

A (binary) symbol code *C*⁶ for an ensemble *X*

 $A_X = \{a, b, c, d\}; P_X = \{1/2, 1/4, 1/8, 1/8\}$ H(X) = 1.75 bits

a_i	$c(a_i)$	p_i	$h(p_i)$	l_i
a	0	$1/_{2}$	1.0	1
b	01	$1/_{4}$	2.0	2
С	011	1/8	3.0	3
d	111	1/8	3.0	3

 C_6 :

 C_3 :

a_i	$c(a_i)$	p_i	$h(p_i)$	l_i
a	0	1/2	1.0	1
b	10	$1/_{4}$	2.0	2
С	110	1/8	3.0	3
d	111	1/8	3.0	3

the codewords of C_6 are the reverse of C_3



What limit is imposed by unique decodeability?



What limit is imposed by unique decodeability?

- Given a list of positive integers $\{l_i\}$, does there exist a uniquely decodeable code with those integers as its codeword lengths?
- We have observed that if we take a code such as {00, 01, 10, 11},
 - shorten one of its codewords, for example $00 \rightarrow 0$,
 - then we can retain unique decodeability only if we lengthen other codewords !
- Thus there seems to be a **constrained budget** that we can spend on codewords, with **shorter**

codewords being more expensive.

- In general for a code with l (constant bits) we have 2^l possible codewords.
- So the "cost" of a codeword of length *l* is 1/2^{*l*}



Kraft inequality

Kraft inequality. For any uniquely decodeable code C(X) over the binary alphabet $\{0, 1\}$,

the codeword lengths must satisfy:

 $\frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} + \frac{1}{2^3} = \frac{1}{2} + \frac{1}{4} + \frac{1}{4} = 1$

$$\sum_{i=1}^{I} 2^{-l_i} \le 1 \quad \text{where} \quad I = |A_X|$$

Completeness. If a uniquely dcodeable codes **satisfies the Kraft inequality with equality** then the code it is called **complete code**.

α	
	٠
$\cup 3$	•

a_i	$c(a_i)$	p_i	$h(p_i)$	l_i
a	0	1/2	1.0	1
b	10	$1/_{4}$	2.0	2
С	110	1/8	3.0	3
d	111	1/8	3.0	3



The symbol coding budget

The "cost" of 2^{-*l*} of each codeword with length *l* is indicated by the size of the box is written

in.

$\begin{array}{ccc} & & & & & & & & & & & & & & & & & &$					
00 000 0001 0001 0001 0001 0001 0001 000 00 0			000	0000	
000 000 0010 0010 0010 0010 0010 0010 0010 0100 0100 100 0010 0100 010 100 0001 0010 0010 100 0001 010 100 100 1001 001 001 010 1001 001 010 101 1001 101 011 111 1110 1111 1111 111		00	000	0001	<u>ب</u>
$\begin{array}{c cccccc} & & & & & & & & & & & & & & & & $		00	001	0010	ge j
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	0		001	0011	nd
1 10 010 010 010 010 010 010 010 010 01	0		010	0100	p p
100 0110 0110 0110 0110 0111 000 0111 000 0111 000 0111 000 0111 000 0100 010		01	010	0101	pde
International Interna International International<		01	011	0110	C S
1 10 1000 1000 1000 1000 1000 1000 100			011	0111	loc
$1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\$		10	100	1000	ml
10 101 1010 1010 1011 1011 1011 1011 1			100	1001	Sy
$1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\$			101	1010	tal
$ \begin{array}{c c} 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\ 1 \\$	1		101	1011	to
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	1		110	1100	he
111 1110 111 1111		11	110	1101	E
			111	1110	
			111	1111	

The **total budget available** when making a uniquely decodeable code is 1.

$$\sum_{i=1}^{I} 2^{-l_i} \le 1$$

If the cost of the codewords that you take exceeds the budget then your code will not be uniquely decodeable.



The symbol coding budget - Code C₀

Code C_0





The symbol coding budget - Code C₃

Code C_3





The symbol coding budget - Code C₄

Code C_4

		000	0000
	00	000	0001
		001	0010
0		001	0011
0		010	0100
	01	010	0101
		011	0110
		011	0111
		100	1000
	10	100	1001
		101	1010
1		101	1011
		110	1100
	11	110	1101
		111	1110
		111	1111

	$\mathbf{\nabla}^{I}$	$-l_i$	1	000	0000
		; '=	•]	000	0001
	i=1 0	0		001	0010
				001	0011
				010	0100
	0	1		010	0101
				011	0110
					0111
				100	1000
	1	0		100	1001
		C_{4}	C_{5}	101	1010
1		04	0.0		1011
1	a	00	0	110	1100
	b	01	1	110	1101
	C	10	00		1110
	d	11	11	111	1111



The symbol coding budget - Code C₆

Code C_6

		000	0000						
	00	000	0001			$\sum 2^{-l_i}$	=1		
	00	001	0010			<i>i</i> =1			
0		001	0011						
0		010	0100						
	01	010	0101						
		011	0110						
			0111						
		100	1000				<u> </u>		
	10	100	1001				C_6 :		
	10	101	1010	-	<i>Q</i> .;	$c(a_i)$	\mathcal{D}_{i}	$h(p_i)$	Li
1		101	1011	-	\sim_l	$\mathcal{O}(\mathcal{O}_{l})$			
1		110	1100		а	0	$\frac{1}{2}$	1.0	1
	11	110	1101		Ъ	01	1/4	2.0	2
	11	111	1110		С	011	1/8	3.0	3
			1111		d	111	1/8	3.0	3



Kraft inequality and prefix codes

- We want codes that are **uniquely decodeable**.
- Prefix codes are uniquely decodeable, and are easy to decode.
- **For any source there is an optimal symbol code that is also a prefix code.**

- Kraft inequality and prefix codes.
 - Given a set of codeword lengths that satisfy the Kraft inequality, there exists a uniquely

decodeable prefix code with these codeword lengths.





• We wish to minimize the expected length of a code

$$L(C,X) = \sum_{i} p_{i} l_{i}$$

Lower bound on expected length. The expected length L(C, X) of a uniquely decodeable code is bounded below by H(X).

Let define the *implicit* probabilities $q_i \equiv 2^{-l_i} / z$, where $z = \sum_{i} 2^{-l_i}$

So
$$l_i = \log \frac{1}{q_i} - \log z$$

Now, consider the relative entropy between *P* and *Q* and the Gibb's inequality

$$D_{KL}(P \parallel Q) = \sum_{x} P(x) \log \frac{P(x)}{Q(x)} \qquad D_{KL}(P \parallel Q) \ge 0$$

= $\sum_{x} P(x) \log P(x) - \sum_{x} P(x) \log Q(x)$
= $\sum_{x} P(x) \log \frac{1}{Q(x)} - \sum_{x} P(x) \log \frac{1}{P(x)}$



Symbol Codes - 28

Ũ

Let define the *implicit* probabilities
$$q_i \equiv 2^{-l_i} / z$$
, where $z = \sum_{i'} 2^{-l_{i'}}$
So $l_i = \log \frac{1}{q_i} - \log z$

Now, consider the relative entropy between *P* and *Q* and the Gibb's inequality

$$D_{KL}(P || Q) = \sum_{x} P(x) \log \frac{1}{Q(x)} - \sum_{x} P(x) \log \frac{1}{P(x)} \qquad D_{KL}(P || Q) \ge 0$$
$$\sum_{i} p_{i} \log \frac{1}{q_{i}} \ge \sum_{i} p_{i} \log \frac{1}{p_{i}}$$

$$L(C,X) = \sum_{i} p_{i}l_{i} = \sum_{i} p_{i}\log \frac{1}{q_{i}} - \log z$$

$$\geq \sum_{i} p_i \log \frac{1}{p_i} - \log z$$

 $\geq H(X)$

The equality L(C, X) = H(X) is achieved only if the Kraft equality z = 1 is satisfied, and if the codelengths satisfy $l_i = \log(1/p_i)$.



Optimal source codelengths. The expected length is minimized and is equal to *H*(*X*) only if the codelengths are equal to the Shannon information contents:

 $l_i = \log_2(1/p_i)$

Implicit probabilities defined by codelengths. Conversely, any choice of codelengths $\{l_i\}$ implicitly defines a probability distribution $\{q_i\}$,

$$q_i \equiv 2^{-l_i} / z$$

for which those codelengths would be the optimal codelengths.

If the code is complete then z = 1 and the implicit probabilities are given by $q_i = 2^{-l_i}$





How much can we compress?



How much can we compress?

So, we can't compress below the entropy. **How close can we expect to get to the entropy?**

Theorem - Source coding theorem for symbol codes.

For an ensemble X there exists a prefix code C with expected length satisfying

 $H(X) \leq L(C, X) \leq H(X) + 1$

Proof.

Let set the codelengths to integers slightly larger than the optimum

$$l_i = \left\lceil \log_2(\frac{1}{p_i}) \right\rceil$$

• We check that there is a prefix code with these lengths by confirming that the Kraft inequality is satisfied



How much can we compress?

 $H(X) \le L(C, X) \le H(X) + 1$

Proof.

- Let set the codelengths to integers slightly larger than the optimum $l_i = \log_2(\frac{1}{p_i})$
- We check that there is a prefix code with these lengths by confirming that the Kraft inequality is satisfied

$$\sum_{i} 2^{-l_{i}} = \sum_{i} 2^{-\left\lceil \log_{2}(\frac{1}{p_{i}}) \right\rceil} \leq \sum_{i} 2^{-\log_{2}(\frac{1}{p_{i}})} = \sum_{i} p_{i} = 1$$

$$\sum_{i=1}^{I} 2^{-l_{i}} \leq 1$$
Kraft inequality
$$\sum_{i} 2^{-\left\lceil \log_{2}(\frac{1}{p_{i}}) \right\rceil} \leq 1$$
Then, there is a prefix code with these l_{i}

Then we can confirm

$$L(C,X) = \sum_{i} p_{i} \left[\log_{2}(\frac{1}{p_{i}}) \right] \leq \sum_{i} p_{i} (\log_{2}(\frac{1}{p_{i}}) + 1) \leq H(X) + 1$$



The cost of using the wrong codelengths

- If we use a **code whose lengths are not equal to the optimal codelengths**, the average message length will be larger than the entropy.
- Let $\{p_i\}$ the true probabilities;
- If we use a complete code with lengths l_i , they define an implicit probabilities $q_i = 2^{-l_i}$

$$L(C, X) = \sum_{i} p_{i}l_{i} = \sum_{i} p_{i}\log(\frac{1}{q_{i}}) \qquad l_{i} = \log_{2}(\frac{1}{q_{i}})$$

$$= H(X) + \sum_{i} p_{i}\log(\frac{1}{q_{i}}) - H(X) \qquad H(X) = \sum_{i} p_{i}\log(\frac{1}{p_{i}})$$

$$= H(X) + \sum_{i} p_{i}\log(\frac{1}{q_{i}}) - \sum_{i} p_{i}\log(\frac{1}{p_{i}})$$

$$= H(X) + \sum_{i} p_{i}\log(\frac{p_{i}}{q_{i}})$$

$$L(C, X) = H(X) + D_{KL}(\mathbf{p} || \mathbf{q})$$
We conclude the extreme best for each probability of the plane.

 \square *L*(*C*, *X*) exceeds the entropy by the relative entropy $D_{\text{KL}}(\mathbf{p} \parallel \mathbf{q})$



Information Theory

Huffman coding



Optimal source coding with symbol codes

- **Goal**: Mminimize the expected length L(C, X).
- How not to do it?
 - One might try to roughly split the set A_X in two, and continue bisecting the subsets so as to define a binary tree from the root.
 - This construction has the right spirit, as in the weighing problem, but it is not necessarily optimal.
 - It achieves $L(C, X) \le H(X) + 2$.

The Huffman coding algorithm

- It is a simple algorithm for finding an optimal prefix code.
- The ideas is to construct the code backwards starting from the tails of the codewords
- The algorithm builds the binary tree from its leaves.



Huffman coding

The Huffman coding algorithm

- It is a simple algorithm for finding an optimal prefix code.
- The ideas is to construct the code backwards starting from the tails of the codewords

The algorithm builds the binary tree from its leaves

- 1. Take the two least probable symbols in the alphabet. These two symbols will be given the longest codewords, which will have equal length, and differ only in the last digit.
- 2. Combine these two symbols into a single symbol, and repeat.
- Since each step reduces the size of the alphabet by one, this algorithm will have assigned

strings to all the symbols after $|A_X| - 1$ steps.



$$\mathcal{A}_X = \{ a, b, c, d, e \}$$
$$\mathcal{P}_X = \{ 0.25, 0.25, 0.2, 0.15, 0.15 \}.$$





$$\mathcal{A}_X = \{ a, b, c, d, e \}$$
$$\mathcal{P}_X = \{ 0.25, 0.25, 0.2, 0.15, 0.15 \}.$$





$$\mathcal{A}_X = \{ a, b, c, d, e \}$$
$$\mathcal{P}_X = \{ 0.25, 0.25, 0.2, 0.15, 0.15 \}.$$

x step 1 step 2 a 0.25 - 0.25 - 0.25b 0.25 - 0.25 - 0.45c 0.2 - 0.2 - 1d 0.15 - 0.3 - 0.3e 0.15 - 1



$$\mathcal{A}_X = \{ a, b, c, d, e \}$$
$$\mathcal{P}_X = \{ 0.25, 0.25, 0.2, 0.15, 0.15 \}.$$

step 1 step 2 step 3 \mathcal{X} 0 0.25 - 0.25 - 0.250.55a 0 0.45 --0.25-0.250.45b 1 - 0.30.2 - 0.2С $0.15 \xrightarrow{0} 0.3$ ′ 1 d 0.15е



$$\mathcal{A}_X = \{ a, b, c, d, e \}$$
$$\mathcal{P}_X = \{ 0.25, 0.25, 0.2, 0.15, 0.15 \}.$$

$$H(X) = 2.2855$$
 bits



Huffman code for English language







Disadvantages of the Huffman code



Disadvantages of the Huffman code

The Huffman algorithm produces an optimal symbol code for an ensemble, but this is not the end of the story. Both the word **'ensemble'** and the phrase **'symbol code'** need careful attention.

Changing ensemble

- Huffman codes do not handle changing ensemble probabilities with any elegance. One brute-force approach would be to recompute the Huffman code every time the probability over symbols changes.
- Or, run through the entire file in advance and compute a good probability distribution, which will then remain fixed throughout transmission. The code itself must also be communicated in this scenario



The extra bit

- Huffman code thus incurs an overhead of between 0 and 1 bits per symbol.
- If H(X) were large, then this overhead would be an unimportant fractional increase. But for many applications, the entropy may be as low as one bit per symbol, or even smaller, so the overhead L(C, X) H(X) may dominate the encoded file length.
- A traditional patch-up of Huffman codes uses them to **compress blocks of symbols**. The overhead per block is at most 1 bit so the overhead per symbol is at most 1/*N* bits.
 - losing the elegant instantaneous decodeability of simple Huffman coding;
 - having to compute the probabilities of all relevant strings and build the associated Huffman tree



Information Theory

Summary



Summary

Kraft inequality. If a code is uniquely decodable its lengths must satisfy

$$\sum_{i} 2^{-l_i} \le 1$$

For any lengths satisfying the Kraft inequality, there exists a prefix code with those

lengths.

Optimal source code lengths for an ensemble are equal to the Shannon information

contents
$$l_i = \log_2 \frac{1}{p_i}$$
,

and conversely, any choice of code lengths defines implicit probabilities

$$q_i = \frac{2^{-l_i}}{z}$$



Summary

- The relative entropy $D_{KL}(\mathbf{p} \parallel \mathbf{q})$ measures how many bits per symbol are wasted by using a code whose implicit probabilities are q, when the ensemble's true probability distribution is p.
- Source coding theorem for symbol codes. For an ensemble *X*, there exists a prefix code whose expected length satisfies

 $H(X) \leq L(C, X) \leq H(X) + 1$

The Huffman coding algorithm generates an optimal symbol code iteratively. At each

iteration, the two least probable symbols are combined.





Further Reading and Summary







Further Reading

Recommend Readings

- Information Theory, Inference, and Learning Algorithms from David MacKay, 2015, pages 91 - 102.
- Supplemental readings:



What you should know

- What are variable-length symbol codes and what it is the main idea
- What are the implications if a symbol code is lossless?
- What is the meaning of a code being uniquely decodeable?
- What are prefix code?
- What is the relation of prefix codes and Kraft inequality?
- What is the the "cost" of 2-I of each codeword?
- What are the compression limits?
- The Huffman coding algorithm
- What the Disadvantages of the Huffman code?



Further Reading and Summary





