4-1-1 INFORMATION PLEASE



"2 bits, 4 bits, 6 bits a byte!"



- Representing Information as bits
- Number Representations
- Other bits

WHAT IS "INFORMATION"?



information, n. Knowledge communicated or received concerning a particular fact or circumstance.





A Computer Scientist's definition:

Information resolves uncertainty.

Information is simply that which cannot be predicted. The less likely a message is, the more information it conveys.

QUANTIFYING INFORMATION

(Claude Shannon, 1948)



Suppose you're faced with N equally probable choices, and I give you a fact that narrows it down to M choices. Then you've been given:

 $log_2(N/M)$ bits of information

Information is measured in bits (binary digits) = number of 0/1's required to encode choice(s)

Examples:

- Outcome of a coin flip: $log_2(2/1) = 1$ bit
- The roll of one die? $\log_2(6/1) = 2.6$ bits
- Someone tells you that their 7-digit phone number is a palindrome?

 $\log_2(10^7/10^4) = ~9.966$ bits

Comp 411 - Fall 2017



ANOTHER EXAMPLE: SUM OF Z DICE



The average information provided by the sum of 2 dice is: 3.274 $i_{ave} = \sum_{i=2}^{12} \frac{M}{N} \log_2(\frac{N}{M_i}) = -\sum_i p_i \log_2(p_i)$ The average information of a process is called its Entropy. 08/25/2017 Comp 411 - Fall 2017

SHOW ME THE BITS!



- Is there a concrete ENCODING that achieves its information content?
- Can the sum of two dice REALLY be represented using 3.274 bits?
- If so, how?
- The fact is, the average information content is a strict lower-bound on how small of a representation that we can achieve.
- In practice, it is difficult to reach this bound. But, we can come very close.



VARIABLE-LENGTH ENCODING



- Of course we can use differing numbers of "bits" to represent each item of data
- This is particularly useful if all items are not equally likely
- Equally likely items lead to fixed length encodings:

 Ex. Encode a "particular" roll of 5?
 {(1,4),(2,3),(3,2),(4,1)} which are equally likely if we use fair dice
 Entropy = \$\sum_{i=1}^4 p(roll_i | roll = 5) log_2(p(roll_i | roll = 5)) = \$\sum_{i=1}^4 \frac{1}{4} log_2(\frac{1}{4}) = 2\$)
 OO = (1,4), OI = (2,3), IO = (3,2), II = (4,1)

 Back to the original problem. Let's use this encoding:
 2 = IOOII
 3 = OIOI
 4 = OII
 5 = OOI
 6 = III
 7 = IOI
 8 = IIO
 9 = 000
 IO = 1000
 II = 0100
 IZ = 10010

VARIABLE-LENGTH DECODING



2 = 10011	3 = 0101	4 = 011	5 = 001
6 = 111	7 = 101	8 = 110	9 = 000
10 = 1000	11 = 0100	12 = 10010	

 Notice how unlikely rolls are encoded using more bits, whereas likely rolls use fewer bits



2 5 3 6 5 8 3 100110010101110011100101

• Where did this code come from?

HUFFMAN CODING



A simple greedy algorithm for approximating a minimum encoding

- 1. Find the 2 items with the smallest probabilities
- 2. Join them into a new *meta* item with probability of their sum
- 3. Remove the two items and insert the new meta item
- 4. Repeat from step 1 until there is only one item



CONVERTING A TREE TO AN ENCODING



Once the *tree* is constructed, label its edges consistently and follow the paths from the largest *meta* item to each of the real items to find the encoding.

2 = 10011	3 = 0101	4 = 011	5 = 001	6 = 111 7 = 101
8 = 110	9 = 000	10 = 1000	11 = 0100	12 = 10010



Comp 411 - Fall 2017

CODING EFFICIENCY



How does this code compare to the information content?

$$b_{ave} = \frac{1}{36}5 + \frac{2}{36}4 + \frac{3}{36}3 + \frac{4}{36}3 + \frac{5}{36}3 + \frac{6}{36}3 + \frac{5}{36}3 + \frac{4}{36}3 + \frac{3}{36}4 + \frac{2}{36}4 + \frac{1}{36}$$
$$b_{ave} = 3.306$$

Pretty close. Recall that the lower bound was 3.274 bits.

However, an efficient encoding (as defined by having an average code size close to the information content) is not always what we want!

Sometimes a uniform code is easier to deal with.

Sometimes redundancy is a good thing.

ENCODING CONSIDERATIONS



- Encoding schemes that attempt to match the information content of a data stream **remove redundancy**. They are **data compression** techniques.
- Make the information easier to manipulate (fixed-sized encodings)
- However, sometimes our goal in encoding information is increase redundancy, rather than remove it. Why?
- Adding redundancy can make data resilient to noise (error detecting and correcting codes)



-Data compression allows us to store our entire music and video collections in a pocketable device -Data redundancy enables us to store that *same* information *reliably* on a hard drive





INFORMATION ENCODING STANDARDS

- "Encoding" is the process of assigning representations to information
- Choosing an appropriate and efficient encoding is an engineering challenge (and an art)
- Impacts design at many levels
 - Mechanism (devices, # of components used)
 - Efficiency (# bits used)
 - Reliability (tolerate noise)
 - Security (encryption)



FIXED-SIZE CODES



If all choices are equally likely (or we have no reason to expect otherwise), then a fixed-size code is often used. Such a code should use at least enough bits to represent the information content. BCD

	0 - 0000
ex. Decimal digits $10 = \{0,1,2,3,4,5,6,7,8,9\}$	1 - 0001
4-bit BCD (binary coded decimal)	2 - 0010
	3 - 0011
$\log_{2}(10/1) = 3.322 < 4 Dits$	4 - 0100
	5 - 0101
ex. ~84 English characters = {A-Z (26), a-z (26),	0-9 (10), 6 - 0110
punctuation (8) m	7 - 0111
	8 - 1000
+Inancial (5)}	9 _ 1001

7-bit ASCII (American Standard Code for Information Interchange) log_(84/1) = 6.392 < 7 bits

9 - 1001





	0000	0001	0010	0011	0100	0101	0110	0111	1000	1001	1010	1011	1100	1101	1110	1111
000	NUL	SOH	STX	ETX	EOT	ACK	ENQ	BEL	BS	НТ	LF	VT	FF	CR	SO	SI
001	DLE	DC1	DC2	DC3	DC4	NAK	SYN	ETB	CAN	EM	SUB	ESC	FS	GS	RS	US
010		!	"	#	\$	%	&		()	*	+	,	-	•	/
011	0	1	2	3	4	5	6	7	8	9	:	;	<	=	>	?
100	@	Α	В	С	D	E	F	G	Н	I	J	К	L	М	N	0
101	Р	Q	R	S	Т	U	V	W	X	Y	Z]	١]	^	_
110	`	а	b	С	d	е	f	g	h	i	j	k	I	m	n	0
111	р	q	r	S	t	u	v	w	x	У	z	{	I	}	~	DEL

- For letters upper and lower case differ in the 6th "shift" bit IOXXXXX is upper, and IIXXXXX is lower
- Special "control" characters set upper two bits to 00 ex. cntl-g \rightarrow bell, cntl-m \rightarrow carriage return, cntl-[\rightarrow esc
- This is why bytes have 8-bits (ASCII + optional parity). Historically, there
 were computers built with 6-bit bytes, which required a special "shift"
 character to set case.

UNICODE



- ASCII is biased towards western languages. English in particular.
- There are, in fact, many more than 256 characters in common use:

â, ö, ß, ñ, è, ¥, £, 揗, 敇, 횝, カ, ‰, 冫, ж, ぢ

- Unicode is a worldwide standard that supports all languages, special characters, classic, extinct, and arcane.
- Several encoding variants 16-bit (UTF-8)
- Variable length (as determined by first byte)



ENCODING POSITIVE INTEGERS



It is straightforward to encode positive integers as a sequence of bits. Each bit is assigned a weight. Ordered from right to left, these weights are increasing powers of 2. The value of an n-bit number encoded in this fashion is given by the following formula:

n-1	211	2 ¹⁰	2 ⁹	28	2	2	6 2	5 2	24	2 ³	2 ²	2 ¹	2 ⁰
$v = \sum_{i=0}^{n} 2^i b_i$	0	1	1	1	1	1	1	0	0	0	0	1	
								2 ⁰ :	_		1		
						-	F	2 ⁵ :	=	3	32		
						-	F	2 ⁶ :	=	6	54		
						4	F	2 ⁷ :	=	12	28		
} { 						-	F	2 ⁸ :	=	2	56		
						4	F	2 ⁹ :	=	5	12		
1. "						-	F ;	2 ¹⁰ :	=	<u>10</u>	<u>24</u>		
- 7										20)16		

Favorite Bits



- You are going to have to get accustomed to working in binary.
 Specifically for Comp 411, but it will be helpful throughout your career as a computer scientist.
- Here are some helpful guides:
 - 1. Memorize the first 10 powers of 2

2 ⁰ = 1	2 ¹ = 2	2 ² = 4	2 ³ = 8	2 ⁴ = 16
2 ⁵ = 32	2 ⁶ = 64	2 ⁷ = 128	2 ⁸ = 256	2 ⁹ = 512

2. Memorize the prefixes for powers of 2 that are multiples of 10

 2^{10} = Kilo (1024) 2^{40} = Tera (1024⁴) 2^{20} = Mega (1024*1024) 2^{50} = Peta (1024⁵) 2^{30} = Giga (1024*1024*1024) 2^{60} = Exa (1024⁶)

TRICKS WITH BITS



- The first thing that you'll do a lot of is cluster groups of contiguous bits.
- Using the binary powers that are multiples of 10 we can do the most basic clustering.
 - 1. When you convert a binary number to decimal, first break it down from the right into clusters of 10 bits.
 - 2. Then compute the value of the leftmost remaining bits (1)
 - 3. Find the appropriate prefix (GIGA)
 - 4. Often this is sufficient (might need to round up)

A "Giga" something or other

OTHER HELPFUL CLUSTERINGS



Oftentimes we will find it convenient to cluster groups of bits together for a more compact written representation. Clustering by 3 bits is called **Octal**, and it is often indicated with a leading zero, **O**. Octal is not that common today.

$v = \sum_{i=0}^{n-1} \xi_{i=0}^{n-1}$	$B^i d_i$	2 ¹¹ 2 ¹⁰ 2 ⁵	^{2⁸ 2⁷ 2⁶ 1111}	2 ⁵ 2 ⁴ 2 ³	2 ² 2 ¹ 2 ⁰	= 2000 ₁₀
	03720	° 3	7	Ý 2	Ŏ	
(Seems natural)	Octal – base	8				
To me!	000 - 0		0*8	3 ⁰ =	0	
p 203	010 - 2		+ 2*8	3 ¹ =	16	
	011 - 3		+ 7*8	$^{2} =$	448	
S-Mert	100 - 4		+ 3*8	3 ³ =	1536	
	101 - 5 110 - 6 111 - 7			2	00010	

ONE MORE CLUSTERING



Clusters of 4 bits are used most frequently. This representation is called hexadecimal. The **hexadecimal** digits include 0-9, and A-F, and each digit position represents a power of 16. Commonly indicated with a leading "Ox".

$v = \sum_{i=0}^{n-1} 16^i d_i$		2 ¹¹ 2 ¹⁰ 2 ⁹ 2 ⁸ 2	⁷ 2 ⁶ 2 ⁵ 2 ⁴ 2 1010	³ 2 ² 2 ¹ 2 ⁰	= 2000 ₁₀
0x7	'dO	7	à	ò	
Hexadecima	l - base 16	5	5	Ū	
0000 - 0 0001 - 1 0010 - 2 0011 - 3 0100 - 4 0101 - 5	1000 - 8 1001 - 9 1010 - a 1011 - b 1100 - c 1101 - d	0 + 13 + 7	*16 ⁰ = 3*16 ¹ = *16 ² =	0 208 1792 2000	0
0110 - 6 0111 - 7	1110 - е 1111 - f			1	0

SUMMARY AND NEXT TIME

- Information is all about bits, Entropy
- Using bits to encode things
 - · Efficient variable-length encodings
 - Redundancy
- Next: more about encoding numbers
 - Signed Numbers (there is a choice)
 - · Non-integers (Fractions and Fixed-point)
 - · Floating point numbers
- Encoding other things ...

