

7

Compression

16 March 2020

Sebastian Wild

Outline

7 Compression

- 7.1 Context
- 7.2 Character Encodings
- 7.3 Huffman Codes
- 7.4 Run-Length Encoding
- 7.5 Lempel-Ziv-Welch
- 7.6 Move-to-Front Transformation
- 7.7 Burrows-Wheeler Transform

7.1 Context

Overview

- ▶ Unit 4–6: How to *work* with strings
 - ▶ finding substrings
 - ▶ finding approximate matches
 - ▶ finding repeated parts
 - ▶ ...
- ▶ Unit 7–8: How to *store* strings
 - ▶ computer memory: must be binary
 - ▶ how to compress strings (save space)
 - (▶ how to robustly transmit over noisy channels \rightsquigarrow Unit 8)

Terminology

- ▶ **source text:** string $S \in \Sigma_S^*$ to be stored / transmitted
 Σ_S is some alphabet
- ▶ **coded text:** encoded data $C \in \Sigma_C^*$ that is actually stored / transmitted
usually use $\Sigma_C = \{0, 1\}$
- ▶ **encoding:** algorithm mapping source texts to coded texts $S \mapsto C$
- ▶ **decoding:** algorithm mapping coded texts back to original source text $C \mapsto S$

What is a good encoding scheme?

- ▶ Depending on the application, goals can be

- ▶ efficiency of encoding/decoding
- ▶ resilience to errors/noise in transmission
- ▶ security (encryption)
- ▶ integrity (detect modifications made by third parties)
- ▶ size

) not here

- ▶ Focus in this unit: size of coded text $|C|$

Encoding schemes that (try to) minimize the size of coded texts perform *data compression*.

- ▶ We will measure the *compression ratio*: $\frac{|C| \cdot \lg |\Sigma_C|}{|S| \cdot \lg |\Sigma_S|} \stackrel{\Sigma_C=\{0,1\}}{=} \frac{|C|}{|S| \cdot \lg |\Sigma_S|}$
 - < 1 means successful compression
 - = 1 means no compression
 - > 1 means “compression” made it bigger!? (yes, that happens ...)

Types of Data Compression

► Logical vs. Physical

- **Logical Compression** uses meaning of data
 - ↪ only applies to a certain domain, e. g., sound recordings
- **Physical Compression** only knows the (physical) **bits** in the data, not the meaning behind them

► Lossy vs. Lossless

- **lossy compression** can only decode approximately;
the exact source text S is lost
- **lossless compression** always decodes S exactly
- For media files, lossy, logical compression is useful (e. g. JPEG, MPEG)
- We will concentrate on physical, lossless compression algorithms.
These techniques can be used for any application.

What makes data compressible?

- ▶ Physical, lossless compression methods mainly exploit two types of redundancies in source texts:

- 1. uneven character frequencies**

some characters occur more often than others → Part I

- 2. repetitive texts**

different parts in the text are (almost) identical → Part II

What makes data compressible?

- Physical, lossless compression methods mainly exploit two types of redundancies in source texts:

1. **uneven character frequencies**

some characters occur more often than others → Part I

2. **repetitive texts**

different parts in the text are (almost) identical → Part II



There is no such thing as a free lunch!

Not *everything* is compressible (→ tutorials)

~> focus on versatile methods that often work

Part I

Exploiting character frequencies

7.2 Character Encodings

Character encodings

- ▶ Simplest form of encoding: Encode each source character individually

↪ encoding function $\underline{E} : \Sigma_S \rightarrow \Sigma_C^*$

- ▶ typically, $|\Sigma_S| \gg |\Sigma_C|$, so need several bits per character
- ▶ for $c \in \Sigma_S$, we call $\underline{E(c)}$ the codeword of c
- ▶ fixed-length code: $|E(c)|$ is the same for all $c \in \Sigma_S$
- ▶ variable-length code: not all codewords of same length

Fixed-length codes

- ▶ fixed-length codes are the simplest type of character encodings
- ▶ Example: ASCII (American Standard Code for Information Interchange, 1963)

| | | | | | | | |
|--------------------|-------------|------------|-----------|------------------|-----------|-----------|-------------|
| 0000000 NUL | 0010000 DLE | 0100000 | 0110000 0 | 1000000 @ | 1010000 P | 1100000 ' | 1110000 p |
| 0000001 SOH | 0010001 DC1 | 0100001 ! | 0110001 1 | 1000001 A | 1010001 Q | 1100001 a | 1110001 q |
| 0000010 STX | 0010010 DC2 | 0100010 " | 0110010 2 | 1000010 B | 1010010 R | 1100010 b | 1110010 r |
| 0000011 ETX | 0010011 DC3 | 0100011 # | 0110011 3 | 1000011 C | 1010011 S | 1100011 c | 1110011 s |
| 0000100 EOT | 0010100 DC4 | 0100100 \$ | 0110100 4 | 1000100 D | 1010100 T | 1100100 d | 1110100 t |
| 0000101 ENQ | 0010101 NAK | 0100101 % | 0110101 5 | 1000101 <u>E</u> | 1010101 U | 1100101 e | 1110101 u |
| 0000110 <u>ACK</u> | 0010110 SYN | 0100110 & | 0110110 6 | 1000110 F | 1010110 V | 1100110 f | 1110110 v |
| 0000111 BEL | 0010111 ETB | 0100111 ' | 0110111 7 | 1000111 G | 1010111 W | 1100111 g | 1110111 w |
| 0001000 BS | 0011000 CAN | 0101000 (| 0111000 8 | 1001000 H | 1011000 X | 1101000 h | 1111000 x |
| 0001001 HT | 0011001 EM | 0101001) | 0111001 9 | 1001001 I | 1011001 Y | 1101001 i | 1111001 y |
| 0001010 LF | 0011010 SUB | 0101010 * | 0111010 : | 1001010 J | 1011010 Z | 1101010 j | 1111010 z |
| 0001011 VT | 0011011 ESC | 0101011 + | 0111011 ; | 1001011 K | 1011011 [| 1101011 k | 1111011 { |
| 0001100 FF | 0011100 FS | 0101100 , | 0111100 < | 1001100 L | 1011100 \ | 1101100 l | 1111100 |
| 0001101 CR | 0011101 GS | 0101101 - | 0111101 = | 1001101 M | 1011101] | 1101101 m | 1111101 } |
| 0001110 SO | 0011110 RS | 0101110 . | 0111110 > | 1001110 N | 1011110 ^ | 1101110 n | 1111110 ~ |
| 0001111 SI | 0011111 US | 0101111 / | 0111111 ? | 1001111 O | 1011111 _ | 1101111 o | 1111111 DEL |

- ▶ 7 bit per character
- ▶ just enough for English letters and a few symbols (plus control characters)

Fixed-length codes – Discussion



Encoding & Decoding as fast as it gets



Unless all characters equally likely, it wastes a lot of space



inflexible (how to support adding a new character?)

Variable-length codes

- ▶ to gain more flexibility, have to allow different lengths for codewords
- ▶ actually an old idea: **Morse Code**

International Morse Code

1. The length of a dot is one unit.
2. A dash is three units.
3. The space between parts of the same letter is one unit.
4. The space between letters is three units.
5. The space between words is seven units.

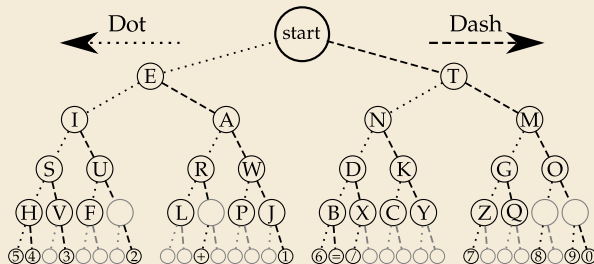
encoding

| | | | |
|---|-----------|---|-------------|
| A | • — | U | • • — |
| B | • • • • | V | • • • — |
| C | • • • • • | W | • • — — |
| D | • • • • • | X | • • • • — |
| E | • | Y | • • • • — • |
| F | • • • • • | Z | • • — — • • |
| G | • • • • • | | |
| H | • • • • • | | |
| I | • • | | |
| J | • — — — — | | |
| K | • • — — — | | |
| L | • • • — — | | |
| M | • • — — | | |
| N | • • • • | | |
| O | • • — — — | | |
| P | • • • — — | | |
| Q | • • — — • | | |
| R | • • • • • | | |
| S | • • • • • | | |
| T | • | | |

https://commons.wikimedia.org/wiki/File:International_Morse_Code.svg

tree

decoding



<https://commons.wikimedia.org/wiki/File:Morse-code-tree.svg>

Variable-length codes – UTF-8

- ▶ Modern example: UTF-8 encoding of Unicode:

default encoding for text-files, XML, HTML since 2009

- ▶ Encodes any Unicode character (137 994 as of May 2019, and counting)
- ▶ uses 1–4 bytes (codeword lengths: 8, 16, 24, or 32 bits)
- ▶ Every ASCII character is encoded in 1 byte with leading bit 0, followed by the 7 bits for ASCII
- ▶ Non-ASCII characters start with 1–4 1s indicating the total number of bytes, followed by a 0 and 3–5 bits.

The remaining bytes each start with 10 followed by 6 bits.

| Char. number range (hexadecimal) | UTF-8 octet sequence (binary) |
|-------------------------------------|-------------------------------------|
| 0000 0000-0000 007F | <u>0xxxxxxx</u> |
| 0000 0080-0000 07FF | 110xxxxx 10xxxxxx |
| 0000 0800-0000 FFFF | 1110xxxx 10xxxxxx 10xxxxxx |
| 0001 0000-0010 FFFF | 11110xxx 10xxxxxx 10xxxxxx 10xxxxxx |



For English text, most characters use only 8 bit,
but we can include any Unicode character, as well.

Pitfall in variable-length codes

- Suppose we have the following code:

| c | a | n | b | s |
|--------|---|----|-----|-----|
| $E(c)$ | 0 | 10 | 110 | 100 |
- Happily encode text $S = \underline{\text{banana}}$ with the coded text $C = \underline{1100} \underline{100} \underline{100}$

b
a
n
a
n
a

Pitfall in variable-length codes

- Suppose we have the following code:

| c | a | n | b | s |
|--------|---|----|-----|-----|
| $E(c)$ | 0 | 10 | 110 | 100 |
- Happily encode text $S = \text{banana}$ with the coded text $C = \underline{1100}\underline{100}\underline{100}$

b
a
n
a
n
a

⚡ C = 1100100100 decodes **both** to banana and to bass: $\frac{1100}{b} \frac{100100}{a \quad s \quad s}$

→ not a valid code ... (cannot tolerate ambiguity)

but how should we have known?

Pitfall in variable-length codes

- ▶ Suppose we have the following code:

| c | a | n | b | s |
|--------|---|----|-----|-----|
| $E(c)$ | 0 | 10 | 110 | 100 |
- ▶ Happily encode text $S = \text{banana}$ with the coded text $C = \underline{1100}\underline{100}\underline{100}$

b
a
n
a
n
a

⚡ $C = 1100100100$ decodes **both** to banana and to bass: $\underline{1100}\underline{100}100$

b
a
s
s

↪ not a valid code ... (cannot tolerate ambiguity)

but how should we have known?



$E(n) = 10$ is a (proper) **prefix** of $E(s) = 100$

101

↪ Leaves decoding wondering whether to stop after reading 10 or continue

↪ Require a prefix-free code: No codeword is a prefix of another.

prefix-free \implies instantaneously decodable

01001
 └────────┘
 codeword

Code tries

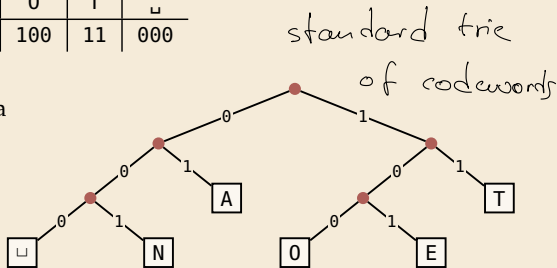
- From now on only consider prefix-free codes E :
 $E(c)$ is not a prefix of $E(c')$ for any $c, c' \in \Sigma_S$.

► Example:

| c | A | E | N | O | T | \sqcup |
|--------|----|-----|-----|-----|----|----------|
| $E(c)$ | 01 | 101 | 001 | 100 | 11 | 000 |

Any prefix-free code corresponds to a
(code) trie (trie of codewords)
with characters of Σ_S at **leaves**.

no need for end-of-string symbols \$ here
(already prefix-free!)



- Encode AN \sqcup ANT 010001000...
- Decode 111000001010111 TO \sqcup

Code tries

- From now on only consider prefix-free codes E :

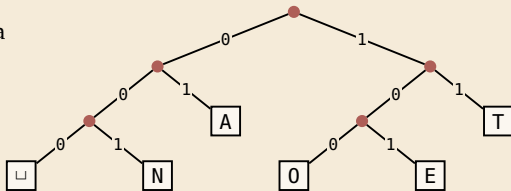
$E(c)$ is not a prefix of $E(c')$ for any $c, c' \in \Sigma_S$.

- **Example:**

| c | A | E | N | O | T | \sqcup |
|--------|----|-----|-----|-----|----|----------|
| $E(c)$ | 01 | 101 | 001 | 100 | 11 | 000 |

Any prefix-free code corresponds to a
(code) trie (trie of codewords)
with characters of Σ_S at **leaves**.

no need for end-of-string symbols $\$$ here
(already prefix-free!)



- Encode $AN_{\sqcup}ANT \rightarrow 010010000100111$
- Decode $1110000001010111 \rightarrow T0_{\sqcup}EAT$

Who decodes the decoder?

- ▶ Depending on the application, we have to **store/transmit** the used code!
- ▶ We distinguish:
 - ▶ fixed coding: code agreed upon in advance, not transmitted (e. g., Morse, UTF-8)
 - ▶ **static coding**: code depends on message, but stays same for entire message; it must be transmitted (e. g., Huffman codes → next)
 - ▶ **adaptive coding**: code depends on message and changes during encoding; implicitly stored withing the message (e. g., LZW → below)

we have
to transmit
code

7.3 Huffman Codes

Character frequencies

- **Goal:** Find character encoding that produces short coded text
- Convention here: fix $\Sigma_C = \{0, 1\}$ (binary codes), abbreviate $\Sigma = \Sigma_S$,
- **Observation:** Some letters occur more often than others.

Typical English prose:

| | | | | | | | | |
|---|--------|----------|---|-------|----|---|-------|---|
| e | 12.70% | ████████ | d | 4.25% | ██ | p | 1.93% | █ |
| t | 9.06% | ██████ | l | 4.03% | ██ | b | 1.49% | █ |
| a | 8.17% | ██████ | c | 2.78% | █ | v | 0.98% | █ |
| o | 7.51% | ██████ | u | 2.76% | █ | k | 0.77% | █ |
| i | 6.97% | ██████ | m | 2.41% | █ | j | 0.15% | |
| n | 6.75% | ██████ | w | 2.36% | █ | x | 0.15% | |
| s | 6.33% | ██████ | f | 2.23% | █ | q | 0.10% | |
| h | 6.09% | ██████ | g | 2.02% | █ | z | 0.07% | |
| r | 5.99% | ██████ | y | 1.97% | █ | | | |

~> Want shorter codes for more frequent characters!

Huffman coding

e. g. frequencies / probabilities

- ▶ **Given:** Σ and weights $w : \Sigma \rightarrow \mathbb{R}_{\geq 0}$
- ▶ **Goal:** prefix-free code E (= code trie) for Σ that minimizes coded text length
i. e., a code trie minimizing $\sum_{c \in \Sigma} w(c) \cdot |E(c)|$

Huffman coding

e. g. frequencies / probabilities

- ▶ **Given:** Σ and weights $w : \Sigma \rightarrow \mathbb{R}_{\geq 0}$
- ▶ **Goal:** prefix-free code E (= code trie) for Σ that minimizes coded text length

i. e., a code trie minimizing $\sum_{c \in \Sigma} w(c) \cdot |E(c)|$

- ▶ If we use $w(c) = \text{\#occurrences of } c \text{ in } S$,
this is the character encoding with smallest possible $|C|$

↪ best possible character-wise encoding

- ▶ Quite ambitious! *Is this efficiently possible?*

Huffman's algorithm

- ▶ Actually, yes! A greedy/myopic approach succeeds here.

Huffman's algorithm:

1. Find two characters a, b with lowest weights.

- ▶ We will encode them with the same prefix, plus one distinguishing bit,

i. e., $E(a) = u0$ and $E(b) = u1$ for a bitstring $u \in \{0, 1\}^*$ (u to be determined)

2. (Conceptually) replace a and b by a single character "ab" \rightarrow Σ decreases by 1
with $w(\text{ab}) = w(a) + w(b)$.

3. Recursively apply Huffman's algorithm on the smaller alphabet.
This in particular determines $u = E(\text{ab})$.

Huffman's algorithm

- ▶ Actually, yes! A greedy/myopic approach succeeds here.

Huffman's algorithm:

1. Find two characters a , b with lowest weights.

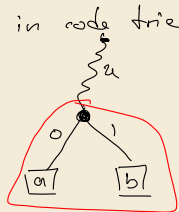
- ▶ We will encode them with the same prefix, plus one distinguishing bit,

i. e., $E(a) = u0$ and $E(b) = u1$ for a bitstring $u \in \{0, 1\}^*$ (u to be determined)

2. (Conceptually) replace a and b by a single character " \boxed{ab} " with $w(\boxed{ab}) = w(a) + w(b)$.

3. Recursively apply Huffman's algorithm on the smaller alphabet. This in particular determines $u = E(\boxed{ab})$.

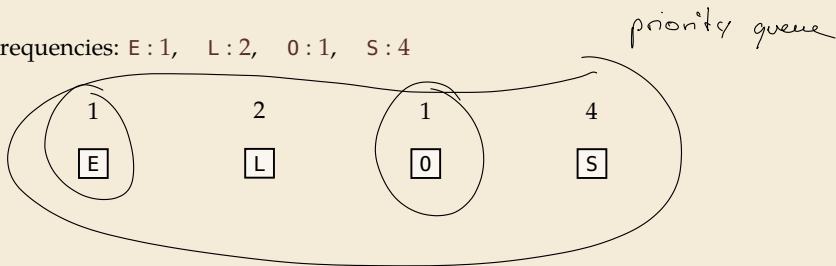
- ▶ efficient implementation using a (min-oriented) priority queue
 - ▶ start by inserting all characters with their weight as key
 - ▶ step 1 uses two deleteMin calls
 - ▶ step 2 inserts a new character with the sum of old weights as key



Huffman's algorithm – Example

► Example text: $S = \underline{\text{LOSSLESS}}$ $\rightsquigarrow \Sigma_S = \{E, L, O, S\}$

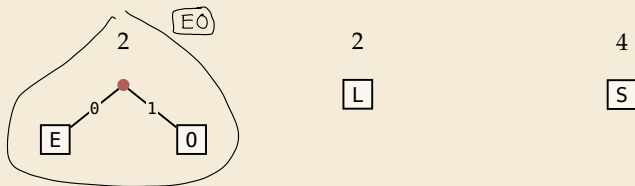
► Character frequencies: E : 1, L : 2, O : 1, S : 4



Huffman's algorithm – Example

► Example text: $S = \text{LOSSLESS}$ $\rightsquigarrow \Sigma_S = \{E, L, O, S\}$

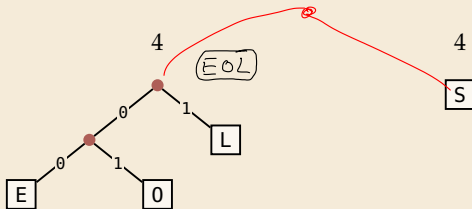
► Character frequencies: $E : 1, \quad L : 2, \quad O : 1, \quad S : 4$



Huffman's algorithm – Example

► Example text: $S = \text{LOSSLESS}$ $\rightsquigarrow \Sigma_S = \{E, L, O, S\}$

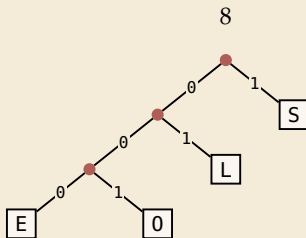
► Character frequencies: E : 1, L : 2, O : 1, S : 4



Huffman's algorithm – Example

► Example text: $S = \text{LOSSLESS}$ $\rightsquigarrow \Sigma_S = \{E, L, O, S\}$

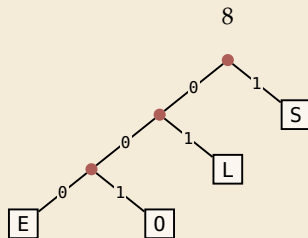
► Character frequencies: E : 1, L : 2, O : 1, S : 4



Huffman's algorithm – Example

► Example text: $S = \text{LOSSLESS}$ $\rightsquigarrow \Sigma_S = \{E, L, O, S\}$

► Character frequencies: E : 1, L : 2, O : 1, S : 4

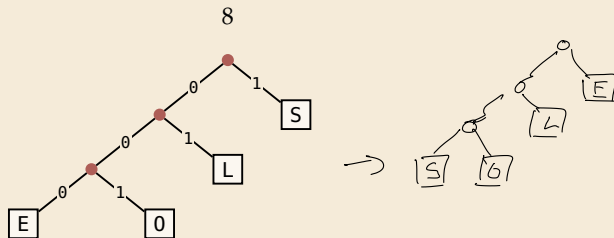


\rightsquigarrow Huffman tree (code trie for Huffman code)

Huffman's algorithm – Example

► Example text: $S = \text{LOSSLESS}$ $\rightsquigarrow \Sigma_S = \{E, L, O, S\}$

► Character frequencies: $E : 1, \quad L : 2, \quad O : 1, \quad S : 4$



\rightsquigarrow *Huffman tree* (code trie for Huffman code)

$\text{LOSSLESS} \rightarrow \underline{01001110100011}$

compression ratio: $\frac{14}{8 \cdot \log 4} = \frac{14}{16} \approx 88\%$

Huffman tree – tie breaking

- ▶ The above procedure is ambiguous:
 - ▶ which characters to choose when weights are equal?
 - ▶ which subtree goes left, which goes right?
- ▶ For COMP 526: always use the following rule:

1. To break ties when selecting the two characters, first use the smallest letter according to the alphabetical order, or the tree containing the smallest alphabetical letter.
2. When combining two trees of different values, place the lower-valued tree on the left (corresponding to a 0-bit).
3. When combining trees of equal value, place the one containing the smallest letter to the left.

Huffman code – Optimality

Theorem 7.1 (Optimality of Huffman's Algorithm)

Given Σ and $w : \Sigma \rightarrow \mathbb{R}_{\geq 0}$, Huffman's Algorithm computes codewords $E : \Sigma \rightarrow \{0, 1\}^*$ with minimal expected codeword length $\ell(E) = \sum_{c \in \Sigma} w(c) \cdot |E(c)|$, among all prefix-free codes for Σ .

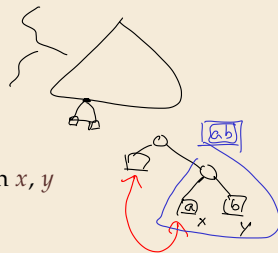
Huffman code – Optimality

Theorem 7.1 (Optimality of Huffman's Algorithm)

Given Σ and $w : \Sigma \rightarrow \mathbb{R}_{\geq 0}$, Huffman's Algorithm computes codewords $E : \Sigma \rightarrow \{0, 1\}^*$ with minimal expected codeword length $\ell(E) = \sum_{c \in \Sigma} w(c) \cdot |E(c)|$, among all prefix-free codes for Σ .

Proof sketch: by induction over $\sigma = |\Sigma|$

- ▶ Given any optimal prefix-free code E^* (as its code trie).
 - ▶ code trie $\rightsquigarrow \exists$ two sibling leaves x, y at largest depth D
 - ▶ swap characters in leaves to have two lowest-weight characters a, b in x, y (that can only make ℓ smaller, so still optimal)
 - ▶ any optimal code for $\Sigma' = \Sigma \setminus \{a, b\} \cup \{ab\}$ yields optimal code for Σ by replacing leaf ab by internal node with children a and b .
- \rightsquigarrow recursive call yields optimal code for Σ' by inductive hypothesis, so Huffman's algorithm finds optimal code for Σ .



Entropy

Definition 7.2 (Entropy)

Given probabilities p_1, \dots, p_n (for outcomes $1, \dots, n$ of a random variable), the *entropy* of the distribution is defined as

$$\mathcal{H}(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i} \right)$$



Entropy

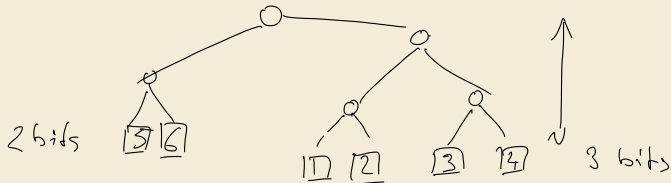
Definition 7.2 (Entropy)

Given probabilities p_1, \dots, p_n (for outcomes $1, \dots, n$ of a random variable), the *entropy* of the distribution is defined as

$$\mathcal{H}(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i} \right)$$

- ▶ entropy is a **measure of information** content of a distribution
 - ▶ more precisely: the expected number of bits (Yes/No questions) required to nail down the random value

↪ would ideally encode value i using $\lg(1/p_i)$ bits
that is not always possible; cannot use 1.5 bits ... but:



fair die

1, 2, ..., 6
 $\frac{1}{6}$... $\frac{1}{6}$

$$\mathcal{H}\left(\frac{1}{6}, \dots, \frac{1}{6}\right) = 6 \cdot \frac{1}{6} \cdot \lg(6) \\ = \lg(6) \approx 2.$$

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

Entropy

Definition 7.2 (Entropy)

Given probabilities p_1, \dots, p_n (for outcomes $1, \dots, n$ of a random variable), the *entropy* of the distribution is defined as

$$\mathcal{H}(p_1, \dots, p_n) = - \sum_{i=1}^n p_i \lg p_i = \sum_{i=1}^n p_i \lg \left(\frac{1}{p_i} \right) \leq \lg(n)$$

- ▶ entropy is a **measure of information** content of a distribution
 - ▶ more precisely: the expected number of bits (Yes/No questions) required to nail down the random value
- ↪ would ideally encode value i using $\lg(1/p_i)$ bits
that is not always possible; cannot use 1.5 bits ... but:

Theorem 7.3 (Entropy bounds for Huffman codes)

For any $\Sigma = \{a_1, \dots, a_\sigma\}$ and $w : \Sigma \rightarrow \mathbb{R}_{\geq 0}$ and its Huffman code E , we have

$$\mathcal{H}\left(\frac{w(a_1)}{W}, \dots, \frac{w(a_\sigma)}{W}\right) \leq \underline{\ell(E)} \leq \mathcal{H}\left(\frac{w(a_1)}{W}, \dots, \frac{w(a_\sigma)}{W}\right) + 1$$

where $W = w(a_1) + \dots + w(a_\sigma)$.

Clicker Question



When is Huffman coding more efficient than a fixed-length encoding?

- ☐ A always
- ☐ B when $\mathcal{H} \approx \lg(\sigma)$
- ☐ C when $\mathcal{H} < \lg(\sigma)$
- ☐ D when $\mathcal{H} < \lg(\sigma) - 1$
- ☐ E when $\mathcal{H} \approx 1$

\mathcal{H} = entropy

$\sigma = |\Sigma|$

pingo.upb.de/622222

Clicker Question



When is Huffman coding more efficient than a fixed-length encoding?

☒ A always ✓

☐ B ~~when $\mathcal{H} \approx \lg(\sigma)$~~

☐ C ~~when $\mathcal{H} < \lg(\sigma)$~~ — $\mathcal{L}(E) \leq \lg \sigma + 1$

☒ D when $\mathcal{H} < \lg(\sigma) - 1$ ✓ $\mathcal{L}(E)$ $\leq \mathcal{H} + 1 < \lceil \lg(\sigma) \rceil$

☐ E ~~when $\mathcal{H} \approx 1$~~

pingo.upb.de/622222

Encoding with Huffman code

- ▶ The overall encoding procedure is as follows:
 - ▶ Pass 1: Count character frequencies in S
 - ▶ Construct Huffman code E (as above)
 - ▶ Store the Huffman code in C (details omitted)
 - ▶ Pass 2: Encode each character in S using E and append result to C
- ▶ Decoding works as follows:
 - ▶ Decode the Huffman code E from C . (details omitted)
 - ▶ Decode S character by character from C using the code trie.
- ▶ Note: Decoding is much simpler/faster!



Huffman coding – Discussion

- ▶ running time complexity: $O(\sigma \log \sigma)$ to construct code
 - ▶ build PQ + $\sigma \cdot (2 \text{ deleteMins and } 1 \text{ insert})$
 - ▶ can do $\Theta(\sigma)$ time when characters already sorted by weight
 - ▶ time for encoding: $O(n + |C|)$
- ▶ many variations in use (tie-breaking rules, estimated frequencies, adaptive encoding, ...)

Huffman coding – Discussion

- ▶ running time complexity: $O(\sigma \log \sigma)$ to construct code
 - ▶ build PQ + $\sigma \cdot (2 \text{ deleteMins and } 1 \text{ insert})$
 - ▶ can do $\Theta(\sigma)$ time when characters already sorted by weight
 - ▶ time for encoding: $O(n + |C|)$
- ▶ many variations in use (tie-breaking rules, estimated frequencies, adaptive encoding, ...)

👍 optimal prefix-free character encoding

👍 very fast decoding

👍 robust encoding flipped bits
 ↓
 local errors only affect 1–2 symbols

👎 needs 2 passes over source text for encoding

- ▶ one-pass variants possible, but more complicated

👎 have to store code alongside with coded text → inflation

Part II

Compressing repetitive texts

7.4 Run-Length Encoding

Run-Length encoding

- ▶ simplest form of repetition: *runs* of characters

[illegible]

same character repeated

- ▶ here: only consider $\Sigma_S = \{0, 1\}$ (work on a binary representation)
 - ▶ can be extended for larger alphabets

Run-Length encoding

- ▶ simplest form of repetition: *runs* of characters

[illegible]

^ same character repeated

- ▶ here: only consider $\Sigma_S = \{0, 1\}$ (work on a binary representation)
 - ▶ can be extended for larger alphabets

→ run-length encoding (RLE):

use runs as phrases: $S = \underbrace{00000}_{\text{run 1}} \underbrace{111}_{\text{run 2}} \underbrace{0000}_{\text{run 3}}$

Run-Length encoding

- ▶ simplest form of repetition: *runs* of characters

[illegible]

^ same character repeated

- ▶ here: only consider $\Sigma_S = \{0, 1\}$ (work on a binary representation)
 - ▶ can be extended for larger alphabets

→ run-length encoding (RLE):

use runs as phrases: $(S = \underbrace{00000}_{\text{00000}} \underbrace{111}_{\text{111}} \underbrace{0000}_{\text{0000}})$

↪ We have to store

- ▶ the first bit of S (either 0 or 1)
- ▶ the length each each run
- ▶ Note: don't have to store bit for later runs since they must alternate.

- Example becomes: 0, 5, 3, 4

Run-Length encoding

- ▶ simplest form of repetition: *runs* of characters

[illegible]

same character repeated

- ▶ here: only consider $\Sigma_S = \{0, 1\}$ (work on a binary representation)
 - ▶ can be extended for larger alphabets

→ run-length encoding (RLE):

use runs as phrases: $S = \underbrace{00000}_{\text{run 1}} \underbrace{111}_{\text{run 2}} \underbrace{0000}_{\text{run 3}}$

→ We have to store

- ▶ the first bit of S (either 0 or 1)
 - ▶ the length each each run
 - ▶ Note: don't have to store bit for later runs since they must alternate.
- ▶ Example becomes: 0, 5, 3, 4

- **Question:** How to encode a run length k in binary? (k can be arbitrarily large!)

Clicker Question



How would you encode a string that can be arbitrarily long?

C string? $c_1 c_2 c_3 c_4 \dots '\backslash 0'$ ^{coding}
\$ ← ∞

Java arrays store length ...
└ (2)

pingo.upb.de/622222

Elias codes

- ▶ Need a prefix-free encoding for $\mathbb{N} = \{1, 2, 3, \dots\}$
 - ▶ must allow arbitrarily large integers
 - ▶ must know when to stop reading

Elias codes

- ▶ Need a prefix-free encoding for $\mathbb{N} = \{1, 2, 3, \dots\}$
 - ▶ must allow arbitrarily large integers
 - ▶ must know when to stop reading
- ▶ But that's simple! Just use unary encoding!
 $7 \mapsto 00000001$ $3 \mapsto 0001$ $0 \mapsto 1$ $30 \mapsto 00000000000000000000000000000001$

Elias codes

- ▶ Need a prefix-free encoding for $\mathbb{N} = \{1, 2, 3, \dots\}$
 - ▶ must allow arbitrarily large integers
 - ▶ must know when to stop reading
- ▶ But that's simple! Just use **unary encoding!**
 $7 \mapsto 00000001$ $3 \mapsto 0001$ $0 \mapsto 1$ $30 \mapsto 00000000000000000000000000000001$
 - 👎 Much too long
 - ▶ (wasn't the whole point of RLE to get rid of long runs??)


Elias codes

- ▶ Need a prefix-free encoding for $\mathbb{N} = \{1, 2, 3, \dots\}$
 - ▶ must allow arbitrarily large integers
 - ▶ must know when to stop reading
- ▶ But that's simple! Just use **unary encoding**!
 $7 \mapsto 00000001$ $3 \mapsto 0001$ $0 \mapsto 1$ $30 \mapsto 00000000000000000000000000000001$
 - 👎 Much too long
 - ▶ (wasn't the whole point of RLE to get rid of long runs??)
- ▶ Refinement: **Elias gamma code**
 - ▶ Store the **length** ℓ of the binary representation in unary
 - ▶ Followed by the binary digits themselves

Elias codes

- ▶ Need a prefix-free encoding for $\mathbb{N} = \{1, 2, 3, \dots\}$
 - ▶ must allow arbitrarily large integers
 - ▶ must know when to stop reading
- ▶ But that's simple! Just use **unary encoding**!

$7 \mapsto 00000001$ $3 \mapsto 0001$ $0 \mapsto 1$ $30 \mapsto 00000000000000000000000000000001$

 Much too long

 - ▶ (wasn't the whole point of RLE to get rid of long runs??)
- ▶ Refinement: ***Elias gamma code***
 - ▶ Store the **length** ℓ of the binary representation in **unary**
 - ▶ Followed by the binary digits themselves
 - ▶ little tricks:
 - ▶ always $\ell \geq 1$, so store $\ell - 1$ instead
 - ▶ binary representation always starts with 1 \rightsquigarrow don't need terminating 1 in unary

\rightsquigarrow Elias gamma code = $\ell - 1$ zeros, followed by binary representation

Examples: $1 \mapsto \underline{1}$, $3 \mapsto \underline{011}$, $5 \mapsto 00101$, $30 \mapsto 000011110$

Clicker Question



Decode the **first** number in Elias gamma code (at the beginning) of the following bitstream:

0001101111011100110.

3 zeros $l=4$

$1101_2 = 13$

pingo.upb.de/622222

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 1$

► Decoding:

$C = 00001101001001010$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$k = 7$

$C = 100111$

► Decoding:

$C = 00001101001001010$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$k = 2$

$C = 100111010$

► Decoding:

$C = 00001101001001010$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$$k = 1$$

$C = 1001110101$

► Decoding:

$$C = 00001101001001010$$
$$S =$$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$k = 20$

$C = 1001110101000010100$

► Decoding:

$C = 00001101001001010$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000000011111111111$

 $k = 11$

$C = 10011101010000101000001011$

► Decoding:

$$C = 00001101001001010$$
$$S =$$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 0$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 0$

$\ell = 3 + 1$

$S =$

Run-length encoding – Examples

► Encoding:

$S = 1111111001000000000000000000000011111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 0$

$\ell = 3 + 1$

$k = 13$

$S = 00000000000000$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 1$

$\ell = 2 + 1$

$k =$

$S = 00000000000000$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 1$

$\ell = 2 + 1$

$k = 4$

$S = 000000000000001111$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 0$

$\ell = 0 + 1$

$k =$

$S = 000000000000001111$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 0000110100100\color{red}1010$

$b = 0$

$\ell = 0 + 1$

$k = 1$

$S = 000000000000001111\color{red}0$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 00001101001001010$

$b = 1$

$\ell = 1 + 1$

$k =$

$S = 0000000000000011110$

Run-length encoding – Examples

► Encoding:

$S = 111111100100000000000000000000001111111111$

$C = 10011101010000101000001011$

Compression ratio: $26/41 \approx 63\%$

► Decoding:

$C = 000011010010010$ **10**

$b = 1$

$\ell = 1 + 1$

$k = 2$

$S = 0000000000000011110$ **11**

Run-length encoding – Discussion

- ▶ extensions to larger alphabets possible (must store next character then)
- ▶ used in some image formats (e. g. TIFF)

Run-length encoding – Discussion

- ▶ extensions to larger alphabets possible (must store next character then)
- ▶ used in some image formats (e. g. TIFF)

👍 fairly simple and fast

👍 can compress n bits to $\Theta(\log n)!$
for extreme case of constant number of runs

👎 negligible compression for many common types of data

- ▶ No compression until run lengths $k \geq 6$
- ▶ **expansion** when run lengths $k = 2$ or 6