# Data Compression = Modeling + Coding

November 3, 2021

Entropy
**Probability Coding**

**symbol codes**
Prefix Code
Relationship to Entropy
Huffman Codes
Combining message

## uniform codes

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

Entropy
**Probability Coding**

**symbol codes**
Prefix Code
Relationship to Entropy
Huffman Codes
Combining message

# Symbol codes

| $i$ | $a_i$ | $p_i$ | | |
|---|---|---|---|---|
| 1 | a | 0.0575 | a | ■ |
| 2 | b | 0.0128 | b | · |
| 3 | c | 0.0263 | c | · |
| 4 | d | 0.0285 | d | · |
| 5 | e | 0.0913 | e | ■ |
| 6 | f | 0.0173 | f | · |
| 7 | g | 0.0133 | g | · |
| 8 | h | 0.0313 | h | · |
| 9 | i | 0.0599 | i | ■ |
| 10 | j | 0.0006 | j | · |
| 11 | k | 0.0084 | k | · |
| 12 | l | 0.0335 | l | · |
| 13 | m | 0.0235 | m | · |
| 14 | n | 0.0596 | n | ■ |
| 15 | o | 0.0689 | o | ■ |
| 16 | p | 0.0192 | p | · |
| 17 | q | 0.0008 | q | · |
| 18 | r | 0.0508 | r | ■ |
| 19 | s | 0.0567 | s | ■ |
| 20 | t | 0.0706 | t | ■ |
| 21 | u | 0.0334 | u | · |
| 22 | v | 0.0069 | v | · |
| 23 | w | 0.0119 | w | · |
| 24 | x | 0.0073 | x | · |
| 25 | y | 0.0164 | y | · |
| 26 | z | 0.0007 | z | · |
| 27 | – | 0.1928 | – | ■ |

Entropy
Probability Coding

symbol codes
Prefix Code
Relationship to Entropy
Huffman Codes
Combining message

## some definitions

### Definition

A *code (source code)* C for a random variable X is a mapping from $\{x_1, x_2, \ldots, x_n\}$ to $\mathcal{D}*$the set of finite-length strings of symbols from an alphabet of length D.

Note:
- $C(x)$ the **codeword** of $x$
- $l(x)$ the length of codeword $C(x)$

### Definition

The expectation length $L(C)$ of a source code $C(X)$ for a random variable X with probability mass function $p(x)$is given by:

$$L(C) = \sum p(x_i)l(x_i) = \sum p_i l_i$$

Entropy
**Probability Coding**

symbol codes
**Prefix Code**
Relationship to Entropy
Huffman Codes
Combining message

## source code examples

### Example

ASCII codes
- 96 printable keyboard characters
- log(96) ??

### Example

Morse code
- special stop symbol

### Definition

A code is called a *prefix code* or *instantaneous code (prefix-free code)* if no codeword is a prefix of any other codeword.

Entropy
Probability Coding

symbol codes
Prefix Code
Relationship to Entropy
Huffman Codes
Combining message

# Kraft-McMillan Inequality.

### Theorem

*For any **uniquely decodable code** C*

$$\sum_{x \in C} 2^{-l(x)} \leq 1,$$

*where $l(x)$ is the length of the codeword. Also, for any set of lengths L such that*

$$\sum_{l \in L} 2^{-l} \leq 1,$$

*there is a prefix code C of the same size such that*

Entropy
**Probability Coding**

symbol codes
Prefix Code
**Relationship to Entropy**
Huffman Codes
Combining message

## Source Coding Theorem

### Theorem

*There exists a variable-length encoding C of an ensemble X such that the average length of an encoded symbol $L(C, X)$ satify:*

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

Entropy
**Probability Coding**

symbol codes
Prefix Code
Relationship to Entropy
**Huffman Codes**
Combining message

## are optimal prefix codes

- Generated from a set of probabilities
- David Huffman developed the algorithm as a student in a 12 class on information theory at MIT in 1950
- The algorithm is the most used component of compression algorithms
- used as the back end of GZIP, JPEG

Entropy
**Probability Coding**

symbol codes
Prefix Code
Relationship to Entropy
**Huffman Codes**
Combining message

## The Huffman algorithm

How it generates the prefix-code tree.

- Start with a forest of trees, one for each message. Each tree contains a single vertex with weight $w_i = p_i$
- Repeat until only a single tree remains
    - Select two trees with the lowest weight roots ($w_i$ and $w_2$)
    - Combine them into a single tree by adding a new root with weight $w_1 + w_2$, and making the two trees its children.

Entropy
**Probability Coding**

symbol codes
Prefix Code
Relationship to Entropy
Huffman Codes
**Combining message**

# Arithmetic coding

- information from the messages is combined to share the same bits.

- asymptotically approach the sum of the self information of the individual messages

- The Shannon information content of an outcome

$$h(x \quad = \quad a_i) = log_2(\frac{1}{P(x = a_i)}) = \log_2 \frac{1}{p_i}$$

- Using a Huffman code, each message has to take at least 1 bit

Entropy
**Probability Coding**

symbol codes
Prefix Code
Relationship to Entropy
Huffman Codes
**Combining message**

# Arithmetic coding

## Example

Entropy
**Probability Coding**

symbol codes
Prefix Code
Relationship to Entropy
Huffman Codes
**Combining message**

# Arithmetic coding

Let $X = \begin{pmatrix} x_1 & \ldots & x_n \\ p_1 & \ldots & p_n \end{pmatrix}$

The **Arithmetic coding algorithm**

- define the accumulated probability

$$f_i = \sum_{j=1}^{i-1} p(j) \text{ with } i = 1, \ldots, n$$

- compute the interval length and size:

$$l_i = \begin{cases} f_i & i = 1 \\ l_{i-1} + f_i * s_{i-1} & 1 < i \leq n \end{cases}$$

$$s_i = \begin{cases} p_i & i = 1 \\ s_{i-1} * p_i & 1 < i \leq n \end{cases}$$