

Dictionary compression

The compression techniques we have seen so far replace individual symbols with a variable length codewords. In dictionary compression, variable length substrings are replaced by short, possibly even fixed length codewords. Compression is achieved by replacing long strings with shorter codewords.

The general scheme is as follows:

- The **dictionary** \mathcal{D} is a collection of strings, often called **phrases**. For completeness, the dictionary includes all single symbols.
- The text T is **parsed** into a sequence of phrases:

$$T = T_1T_2\dots T_z, T_i \in \mathcal{D}.$$

The sequence is called a **parsing** or a **factorization** of T with respect to \mathcal{D} .

- The text is encoded by replacing each phrase T_i with a code that acts as a pointer to the dictionary.

Here is a simple static dictionary encoding for English text:

- The dictionary consists of some set of English words plus individual symbols.
- Compute the frequencies of the words in some corpus of English texts. Compute the frequencies of symbols in the corpus from which the dictionary words have been removed.
- Number the words and symbols in descending order of their frequencies.
- To encode a text, replace each dictionary word and each symbol that does not belong to a word with its corresponding number. Encode the sequence of number using γ coding.

Lempel-Ziv compression

In 1977 and 1978, Abraham Lempel and Jacob Ziv published two adaptive dictionary compression algorithms that soon became to dominate practical text compression. Numerous variants have been published and implemented, and they are still the most commonly used algorithms in general purpose compression tools.

The common feature of the two algorithms and all their variants is that the dictionary consists of substrings of the already processed part of the text. This means that the dictionary adapts to the text.

The two algorithms are known as LZ77 and LZ78, and most related methods can be categorized as a variant of one or the other. The primary difference is the encoding of the phrases:

- LZ77 uses direct pointers to the preceding text.
- LZ78 uses pointers to a separate dictionary.

LF78

The original LZ78 encoding works as follows:

- The dictionary consists of phrases numbered from 0 upwards: $\mathcal{D} = \{Z_0, Z_1, Z_2, \dots\}$. Each new phrase inserted to the dictionary gets the next free number. Initially, the only phrase is the empty string $Z_0 = \epsilon$.
- Suppose we have so far computed the parsing $T_1 \dots T_{j-1}$ for $T[0 \dots i)$ and the next phrase T_j starts at position i . Let $Z_k \in \mathcal{D}$ be the longest phrase in the dictionary that is a prefix of $T[i \dots n - 1)$. The next phrase is $T_j = Z_j = T[i \dots i + |Z_k|] = Z_k t_{i+|Z_k|}$, and it is inserted to the dictionary.
- The phrase T_j is encoded as the pair $\langle k, t_{i+|Z_k|} \rangle$. Using fixed length codes, the pair needs $\lceil \log(j + 1) \rceil + \lceil \log \sigma \rceil$ bits.

Example 2.7: Let $T = \text{badadadabaab}$.

index	0	1	2	3	4	5	6	7
phrase	ϵ	b	a	d	ad	ada	ba	ab
encoding		$\langle 0, \mathbf{b} \rangle$	$\langle 0, \mathbf{a} \rangle$	$\langle 0, \mathbf{d} \rangle$	$\langle 2, \mathbf{d} \rangle$	$\langle 4, \mathbf{a} \rangle$	$\langle 1, \mathbf{a} \rangle$	$\langle 2, \mathbf{b} \rangle$
code length		0+2	1+2	2+2	2+2	3+2	3+2	3+2

LZW

LZW is simple optimization of LZ78 used, e.g., in the unix tool compress.

- Initially, the dictionary \mathcal{D} contains all individual symbols:
 $\mathcal{D} = \{Z_1, \dots, Z_\sigma\}$.
- Suppose the next phrase T_j starts at position i . Let $Z_k \in \mathcal{D}$ be the longest phrase in the dictionary that is a prefix of $T[i..n]$. Now the next text phrase is $T_j = Z_k$ and the phrase added to the dictionary is $Z_{\sigma+j} = T_j t_{i+|T_j|}$.
- The phrase T_j is encoded with the index k requiring $\lceil \log(\sigma + j - 1) \rceil$ bits.

Omitting the symbol codes saves space in practice, even though the index codes can be longer and the phrases can be shorter.

Example 2.8: Let $T = \text{badadadabaab}$.

text phrase											
		b	a	d	ad	ada	ba	a	b		
encoding		2	1	4	6	8	5	1	2		
code length		2	3	3	3	3	4	4	4		
dictionary phrase	a	b	c	d	ba	ad	da	ada	adab	baa	ab
index	1	2	3	4	5	6	7	8	9	10	11

There is a tricky detail in decoding LZW:

- In order to insert the phrase $Z_{\sigma+j} = T_j t_{i+|T_j|}$ into the dictionary, the decoder needs to know the symbol $t_{i+|T_j|}$. For the encoder this is not a problem, but the decoder only knows that $t_{i+|T_j|}$ is the first symbol of the next phrase T_{j+1} . Thus the insertion is delayed until the next phrase has been decoded.
- The encoder inserts $Z_{\sigma+j}$ to the dictionary without a delay and has the option of choosing $T_{j+1} = Z_{\sigma+j}$. If this happens, the decoder is faced with a reference to a phrase that is not yet in the dictionary! However, in this case the decoder knows that the unknown symbol $t_{i+|T_j|}$ is the first symbol of $Z_{\sigma+j}$, which is the same as the first symbol of T_j . Given $t_{i+|T_j|}$ the decoder can set $Z_{\sigma+j} = T_j t_{i+|T_j|}$, insert it into the dictionary, and decode T_{j+1} normally.

The phrase $T_{j+1} = Z_{\sigma+j}$ is problematic because it is [self-referential](#).

Example 2.9: In Example 2.8, the phrase $T_5 = Z_8 = \text{ada}$ is self-referential.