

What are the Ingredients of a Theory of Data Compression?

- A context for the problem.
Eg, what are we trying to compress, and what are we compressing it into?
- A notion of what data compression schemes are *possible*.
A data compression scheme must allow us to *encode* data, and then *decode* it, recovering the original data.
- A measure of *how good* a data compression scheme is.
We will have to look at how good a scheme is *on average*, given some model for the source.

One Danger: If we don't formalize things well, we might eliminate data compression schemes that would have been practical.

What Do We Hope to Get From a Theory of Data Compression?

- Easier ways of telling whether a data compression scheme is possible, and if so, how good it is.
- A theorem that tells us how good a scheme can possibly be — the “theoretical limit”.
- Some help in finding a scheme that approaches this theoretical limit.
- Insight into the nature of the problem, which may help for other problems.

One insight: Compression is limited by the *entropy* of the source, which is a measure of *information content* that has many other uses.

Formalizing the Source of Data

We'll assume that we are trying to compress data from a digital *source* that produces a sequence of symbols, X_1, X_2, X_3, \dots

These *source symbols* come from a finite *source alphabet*, $S = \{s_1, \dots, s_q\}$.

Examples:

$$S = \{A, B, \dots, Z, _ \}$$

$$S = \{0, 1, 2, \dots, 255\}$$

$$S = \{C, G, T, A\}$$

$$S = \{0, 1\}$$

The source alphabet is known to the receiver — who may be us at a later time, for storage applications.

Formalizing What We Compress To

The output of the compression program is a sequence of *code symbols* from a finite *code alphabet*, $T = \{t_1, \dots, t_r\}$.

These symbols are sent through the *channel*, to the receiver. We assume for now that the channel is noise-free — the symbol received is always the symbol that was sent.

We'll almost always assume that $T = \{0, 1\}$, since computer files and digital transmissions are usually binary, but the theory applies to any finite T .

Possible Compression Programs

A compression program (ie, a *code*) defines a mapping of each source symbol to a finite sequence of code symbols (a *codeword*).

For example:

$$S = \{C, G, T, A\}, \quad T = \{0, 1\}$$

$$C \rightarrow 0$$

$$G \rightarrow 10$$

$$T \rightarrow 110$$

$$A \rightarrow 1110$$

We encode a sequence of source symbols by concatenating the codewords obtained by this mapping. For example:

$$CCAT \rightarrow 001110110$$

We require that the mapping be such that we can *decode* this sequence.

Later, we'll see that the above formalization isn't really right...

What Codes are Decodable?

We intend to consider only codes that can be decoded. But what do we mean by that?

This may depend on how the channel behaves at the end of a transmission. Four possibilities:

- The end of the transmission is explicitly marked, say by "\$":

011101101\$

- After the end of the transmission, subsequent symbols all have a single known value, say "0":

011101101000000000...

- After the end of the transmission, subsequent symbols are random garbage:

0111011011100100101...

- There is no end to the transmission.

When Do We Need the Decoding?

Another possible issue is when we require that a decoded symbol be known. Possibilities:

- As soon as the codeword for the symbol has been received.
If this is possible, the code is *instantaneously decodable*.
- With no more than a fixed delay after the codeword for the symbol has been received.
If this is possible, the code is *decodable with bounded delay*.
- Not until the entire message has been received.

Assuming that the end of transmission is explicitly marked, we then require only that the code be *uniquely decodable*.

How Much Difference Does it Make?

We could develop theories of data compression with various definitions of decodability.

Question: How much difference will it make?

Will we find that we can't compress data as much if we insist on using a code that is instantaneously decodable?

Or will we find that a single theory is "robust" — not sensitive to the exact details of the channel and decoding requirements.

Easiest: Assume the end of transmission is explicitly marked; don't require any symbols be decoded until the entire message is received.

Hardest: Require instantaneous decoding. (It then won't matter whether the end of transmission is marked, as far as decoding the symbols that were actually sent is concerned.)

Notation for Sequences & Codes

S and T are the source and code alphabets.

S^* and T^* denote sequences of **zero** or more symbols from the source or code alphabets.

S^+ and T^+ denote sequences of **one** or more symbols from the source or code alphabets.

A code (ie, a compression program), \mathcal{C} , is a mapping $S \rightarrow T^+$.

We can extend this to a mapping $\mathcal{C} : S^* \rightarrow T^*$ using concatenation:

$$\mathcal{C}(s_{i_1}s_{i_2}\cdots s_{i_n}) = \mathcal{C}(s_{i_1})\mathcal{C}(s_{i_2})\cdots\mathcal{C}(s_{i_n})$$

We sometimes also use \mathcal{C} to denote the set of codewords: $\{w_i \mid w_i = \mathcal{C}(s_i) \text{ for some } s_i \in S\}$.

Formalizing Uniquely Decodable and Instantaneous Codes

We can now define a code to be *uniquely decodable* if the mapping $\mathcal{C} : S^* \rightarrow T^*$ is one-to-one — ie, if each $t \in T^*$ corresponds to at most one $s \in S^*$.

A code is obviously not uniquely decodable if two symbols have the same codeword — ie, if $\mathcal{C}(s_i) = \mathcal{C}(s_j)$ for some $i \neq j$ — so we'll usually assume that this isn't the case.

We define a code to be *instantaneously decodable* if any source sequences s and s' in S^+ for which s is not a prefix of s' have encodings $t = \mathcal{C}(s)$ and $t' = \mathcal{C}(s')$ for which t is not a prefix of t' . (Since otherwise, after receiving t , we wouldn't yet know whether the message starts with s or with s' .)

Examples

Examples with $S = \{x, y, z\}$ and $T = \{0, 1\}$:

	Code A	Code B	Code C	Code D
x	10	0	0	0
y	11	10	01	01
z	111	110	011	11

Code A: Not uniquely decodable
Both yyy and zz encode as 111111

Code B: Instantaneously decodable
End of each codeword marked by 0

Code C: Decodable with one-symbol delay
End of codeword marked by *following* 0

Code D: Uniquely decodable, but with unbounded delay:
0111111111111111 decodes as $xzzzzzzz$
0111111111111111 decodes as $yzzzzzz$

How to Check Whether a Code is Uniquely Decodable

The *Sardinas-Patterson Theorem* tells us how to check whether a code is uniquely decodable.

Let \mathcal{C} be the set of codewords.

Define $\mathcal{C}_0 = \mathcal{C}$.

For $n > 0$, define

$$\mathcal{C}_n = \{w \in T^+ \mid uw = v \text{ where } u \in \mathcal{C}, v \in \mathcal{C}_{n-1} \text{ or } u \in \mathcal{C}_{n-1}, v \in \mathcal{C}\}$$

Finally, define

$$\mathcal{C}_\infty = \mathcal{C}_1 \cup \mathcal{C}_2 \cup \mathcal{C}_3 \cup \dots$$

The code \mathcal{C} is uniquely decodable if and only if \mathcal{C} and \mathcal{C}_∞ are disjoint.

Applying This Check to the Examples

Code A: $\mathcal{C} = \mathcal{C}_0 = \{10, 11, 111\}$
 $\mathcal{C}_1 = \{1\}$
 $\mathcal{C}_2 = \{0, 1, 11\}$
 $\Rightarrow 11 \in \mathcal{C} \cap \mathcal{C}_\infty$

Code B: $\mathcal{C} = \mathcal{C}_0 = \{0, 10, 110\}$
 $\mathcal{C}_1 = \emptyset$
 $\Rightarrow \mathcal{C}_\infty = \emptyset$

Code C: $\mathcal{C} = \mathcal{C}_0 = \{0, 01, 011\}$
 $\mathcal{C}_1 = \{1, 11\}$
 $\mathcal{C}_2 = \emptyset$
 $\Rightarrow \mathcal{C}_\infty = \{1, 11\}$,
disjoint from \mathcal{C}

Code D: $\mathcal{C} = \mathcal{C}_0 = \{0, 01, 11\}$
 $\mathcal{C}_1 = \{1\}$
 $\mathcal{C}_2 = \{1\}$
 $\Rightarrow \mathcal{C}_\infty = \{1\}$,
disjoint from \mathcal{C}

How to Check Whether a Code is Instantaneously Decodable

A code is instantaneous if and only if no codeword is a prefix of some other codeword.

Proof:

(\Rightarrow) If codeword $\mathcal{C}(s_i)$ is a prefix of codeword $\mathcal{C}(s_j)$, then the encoding of the sequence $s = s_i$ is obviously a prefix of the encoding of the sequence $s' = s_j$.

(\Leftarrow) If the code is not instantaneous, let $t = \mathcal{C}(s)$ be an encoding that is a prefix of another encoding $t' = \mathcal{C}(s')$, but with s not a prefix of s' , and let s be as short as possible.

The first symbols of s and s' can't be the same, since if they were, we could drop these symbols and get a shorter instance. So these two symbols must be different, but one of their codewords must be a prefix of the other.