# CSC 310: Information Theory

University of Toronto, Fall 2011

Instructor: Radford M. Neal

Week 2

# What's Needed for a Theory of (Lossless) Data Compression?

- A context for the problem.

  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 Might We Hope for 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, $X_1$, $X_2$, …. These will be viewed as *random variables*; some particular values they take on will be denoted by $x_1$, $x_2$, ….

These *source symbols* are from a finite *source alphabet*, $\mathcal{A}_X$.

Examples:

$$\mathcal{A}_X = \{A,\ B,\ \ldots,\ Z,\ \_\}$$
$$\mathcal{A}_X = \{0,\ 1,\ 2,\ \ldots,\ 255\}$$
$$\mathcal{A}_X = \{C,\ G,\ T,\ A\}$$
$$\mathcal{A}_X = \{0,\ 1\}$$

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*, $Z_1, Z_2, \ldots$ from a finite *code alphabet*, $\mathcal{A}_Z$.

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 the code alphabet is $\{0, 1\}$, since computer files and digital transmissions are usually binary, but the theory can easily be generalized to any finite code alphabet.

# 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, suppose our source alphabet is $\mathcal{A}_X = \{C, G, T, A\}$.

One possible code is

$$
\begin{aligned}
C &\rightarrow 0 \\
G &\rightarrow 10 \\
T &\rightarrow 110 \\
A &\rightarrow 1110
\end{aligned}
$$

We encode a sequence of source symbols by concatenating the codewords of each:

$$CCAT \rightarrow 001110110$$

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?

Let's consider only codes that can be decoded. What does that mean?

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":

    0111011010000000000 ...

- 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

$\mathcal{A}_X$ and $\mathcal{A}_Z$ are the source and code alphabets.

$\mathcal{A}_X^+$ and $\mathcal{A}_Z^+$ denote sequences of one or more symbols from the source or code alphabets.

A symbol code, $C$, is a mapping $\mathcal{A}_X \to \mathcal{A}_Z^+$. We use $c(x)$ to denote the codeword $C$ maps $x$ to.

We can use concatenation to extend this to a mapping for the textitextended code, $C^+ : \mathcal{A}_X^+ \to \mathcal{A}_Z^+$:

$$c^+(x_1 x_2 \cdots x_N) \quad = \quad c(x_1)c(x_2)\cdots c(x_N)$$

That is, we code a string of symbols by just stringing together the codes for each symbol.

We sometimes also use $C$ to denote the set of codewords:

$$\{w \mid w = c(a) \text{ for some } a \in \mathcal{A}_X\}$$

# Formalizing Uniquely Decodable and Instantaneous Codes

We can now define a code to be *uniquely decodable* if the mapping $C^+ : \mathcal{A}_X^+ \to \mathcal{A}_Z^+$ is one-to-one. In other words:

For all $\boldsymbol{x}$ and $\boldsymbol{x}'$ in $\mathcal{A}_X^+$, $\boldsymbol{x} \neq \boldsymbol{x}'$ imples $c^+(\boldsymbol{x}) \neq c^+(\boldsymbol{x}')$

A code is obviously not uniquely decodable if two symbols have the same codeword — ie, if $c(a) = c(a')$ for some $a \neq a'$ — so we'll usually assume that this isn't the case.

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

# Examples

Examples with $\mathcal{A}_X = \{a,\ b,\ c\}$ and $\mathcal{A}_Z = \{0, 1\}$:

|   | Code A | Code B | Code C | Code D |
|---|--------|--------|--------|--------|
| $a$ | 10  | 0   | 0   | 0  |
| $b$ | 11  | 10  | 01  | 01 |
| $c$ | 111 | 110 | 011 | 11 |

Code A: Not uniquely decodable

Both *bbb* and *cc* 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 *accccccc*

011111111111111 decodes as *bccccccc*

# 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 $C$ be the set of codewords.

Define $C_0 = C$.

For $n > 0$, define

$$C_n = \{w \in \mathcal{A}_Z^+ \mid uw = v \text{ where } u \in C,\ v \in C_{n-1} \text{ or } u \in C_{n-1},\ v \in C\}$$

Finally, define

$$C_\infty = C_1 \cup C_2 \cup C_3 \cup \cdots$$

The code $C$ is uniquely decodable if and only if $C$ and $C_\infty$ are disjoint.

We won't both much with this theorem, since as we'll see it isn't of much practical use.

# 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 $\mathcal{C}(a)$ is a prefix of codeword $\mathcal{C}(a')$, then the encoding of the sequence $\boldsymbol{x} = a$ is obviously a prefix of the encoding of the sequence $\boldsymbol{x}' = a'$.

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

The first symbols of $\boldsymbol{x}$ and $\boldsymbol{x}'$ 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.

# Existence of Codes With Given Lengths of Codewords

Since we hope to compress data, we would like codes that are uniquely decodable and whose codewords are short.

If we could make all the codewords really short, life would be really easy. Too easy.

Instead, making some codewords short will require that other codewords be long, if the code is to be uniquely decodable.

**Questions:** What sets of codeword lengths are possible? Is the answer to this question different for instantaneous codes than for uniquely decodable codes?

# McMillan's Inequality

There is a uniquely decodable binary code with codewords having lengths $l_1, \ldots, l_I$ if and only if

$$\sum_{i=1}^{I} \frac{1}{2^{l_i}} \leq 1$$

**Examples:**

There is a uniquely decodable binary code with lengths 1, 2, 3, 3, since

$$1/2 + 1/4 + 1/8 + 1/8 \; = \; 1$$

An example of such a code is $\{0, 01, 011, 111\}$.

There is *no* uniquely decodable binary code with lengths 2, 2, 2, 2, 2, since

$$1/4 + 1/4 + 1/4 + 1/4 + 1/4 \; > \; 1$$

# Kraft's Inequality

There is an instantaneous binary code with codewords having lengths $l_1, \ldots, l_I$ if and only if

$$\sum_{i=1}^{I} \frac{1}{2^{l_i}} \leq 1$$

**Examples:**

There is an instantaneous binary code with lengths 1, 2, 3, 3, since

$$1/2 + 1/4 + 1/8 + 1/8 = 1$$

An example of such a code is $\{0, 10, 110, 111\}$.

There is an instantaneous binary code with lengths 2, 2, 2, since

$$1/4 + 1/4 + 1/4 < 1$$

An example of such a code is $\{00, 10, 01\}$.

# Implications for Instantaneous and Uniquely Decodable Codes

Combining Kraft's and McMillan's inequalities, we conclude that there is an instantaneous binary code with lengths $l_1, \ldots, l_I$ if and only if there is a uniquely decodable code with these lengths.

**Implication:** There is probably no practical benefit to using uniquely decodable codes that aren't instantaneous.

# Proving the Two Inequalities

We can prove both Kraft's and McMillan's inequality by proving that for any set of lengths, $l_1, \ldots, l_I$, for binary codewords:

A) If $\displaystyle\sum_{i=1}^{I} 1/2^{l_i} \leq 1$, we can construct an instantaneous code with codewords having these lengths.

B) If $\displaystyle\sum_{i=1}^{I} 1/2^{l_i} > 1$, there is no uniquely decodable code with codewords having these lengths.
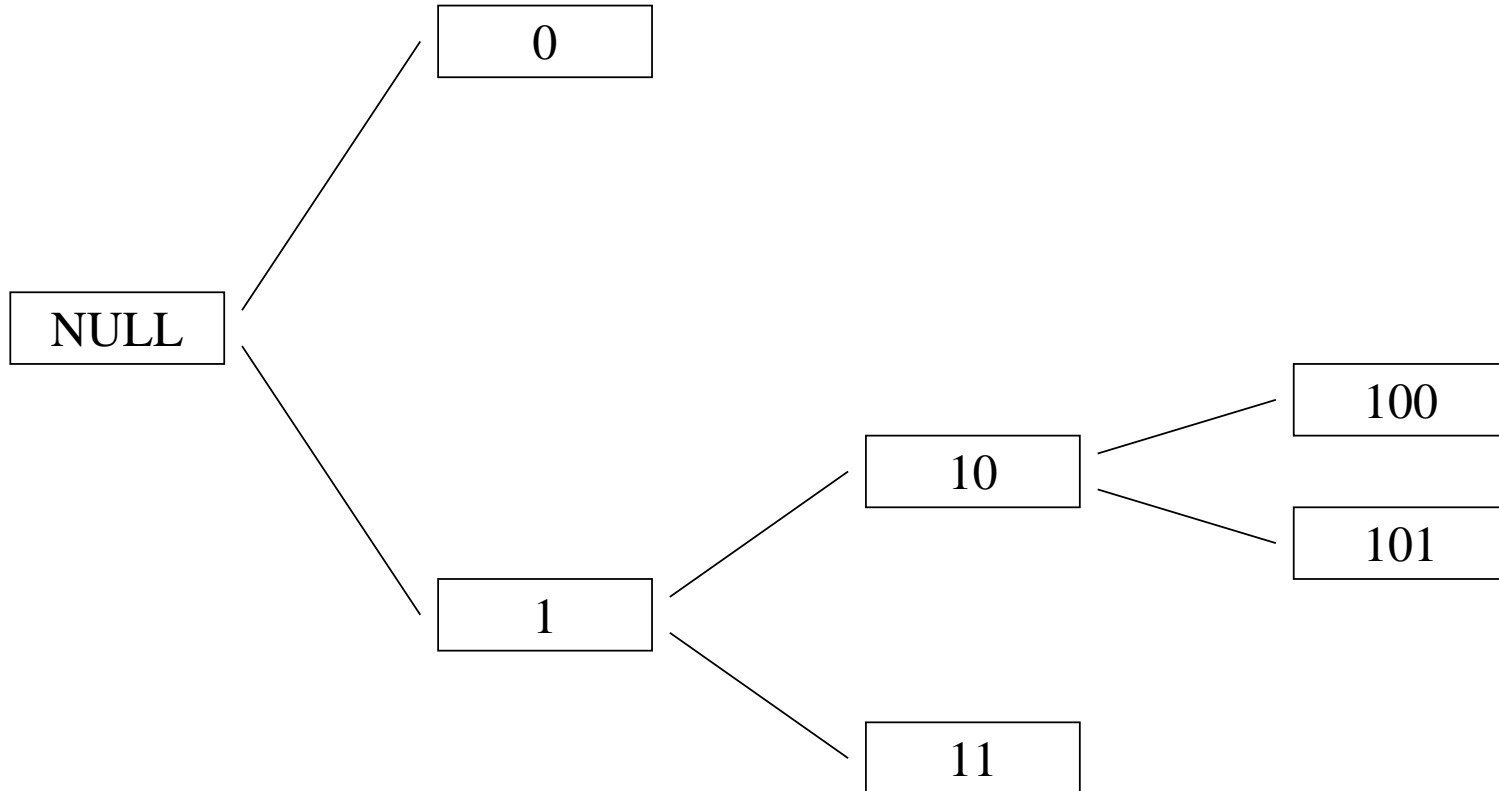
(A) is half of Kraft's inequality. (B) is half of McMillan's inequality.

Since instantaneous codes are uniquely decodable, we also see that (A) gives the other half of McMillan's inequality, and (B) gives the other half of Kraft's inequality.

# Visualizing Prefix Codes as Trees

We can view codewords of an instantaneous (prefix) code as leaves of a tree. The root represents the null string; each branch corresponds to adding another code symbol.
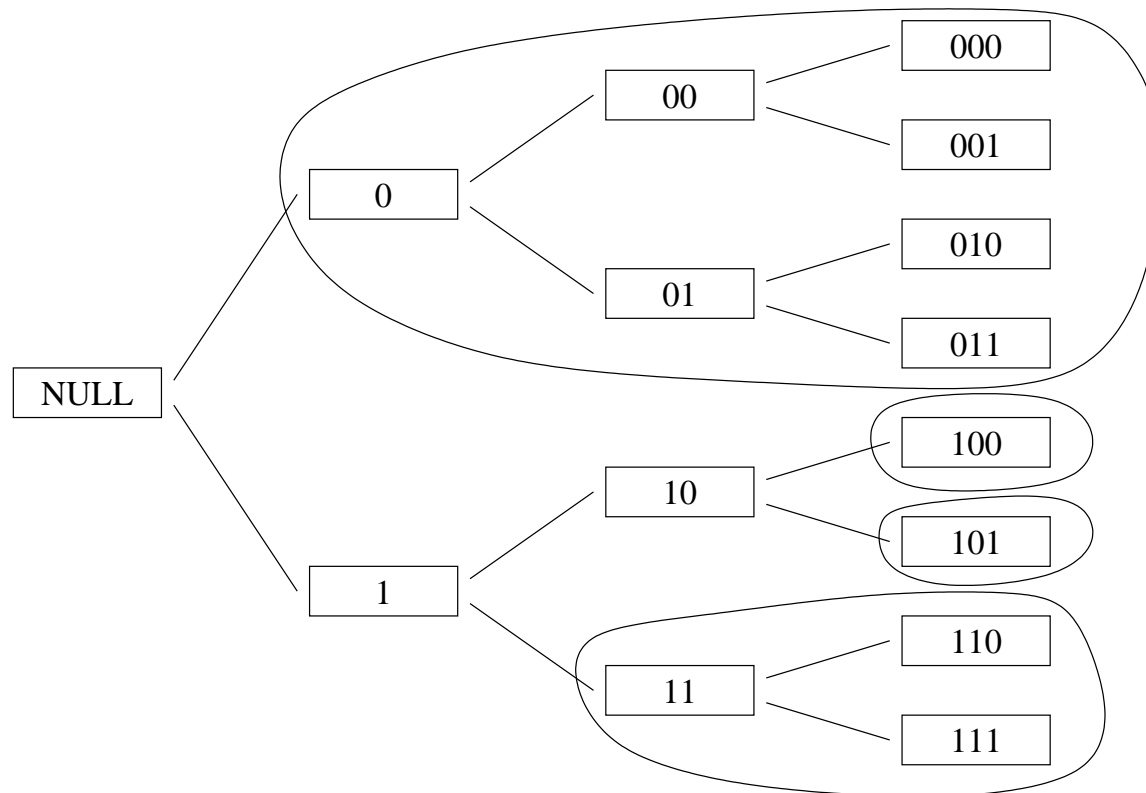
Here is the tree for a code with codewords 0, 11, 100, 101:

# Extending the Tree to Maximum Depth

We can extend the tree to the depth of the longest codeword. Each codeword then corresponds to a subtree.

The extension of the previous tree, with each codeword's subtree circled:



Short codewords occupy more of the tree. For a binary code, the fraction of leaves taken by a codeword of length $l$ is $1/2^l$.

# Constructing Instantaneous Codes
## When the Inequality Holds

Suppose that Kraft's Inequality holds:

$$\sum_{i=1}^{I} \frac{1}{2^{l_i}} \leq 1$$

Order the lengths so $l_1 \leq \cdots \leq l_I$. In the binary tree with depth $l_I$, how can we allocate subtrees to codewords with these lengths?

We go from shortest to longest, $i = 1, \ldots, I$:

1) Pick a node at depth $l_i$ that isn't in a subtree previously used, and let the code for codeword $i$ be the one at that node.

2) Mark all nodes in the subtree headed by the node just picked as being used, and not available to be picked later.

Will there always be a node available in step (1) above?

# Why the Construction Will be Possible

If Kraft's inequality holds, we will always be able to do this.

To begin, there are $2^{l_b}$ nodes at depth $l_b$.

When we pick a node at depth $l_a$, the number of nodes that become unavailable at depth $l_b$ (assumed not less than $l_a$) is $2^{l_b - l_a}$.

When we need to pick a node at depth $l_j$, after having picked earlier nodes at depths $l_i$ (with $i < j$ and $l_i \leq l_j$), the number of nodes left to pick from will be

$$2^{l_j} - \sum_{i=1}^{j-1} 2^{l_j - l_i} \;=\; 2^{l_j}\left[1 - \sum_{i=1}^{j-1} \frac{1}{2^{l_i}}\right] \;>\; 0$$

Since $\displaystyle\sum_{i=1}^{j-1} \frac{1}{2^{l_i}} < \sum_{i=1}^{I} \frac{1}{2^{l_i}} \leq 1$, by assumption.

# Why Uniquely Decodable Codes Must Obey the Inequality

Let $l_1 \leq \cdots \leq l_I$ be the codeword lengths. Define $K = \sum_{i=1}^{I} \frac{1}{2^{l_i}}$.

For any positive integer $n$,

$$K^n = \left[\sum_{i=1}^{I} \frac{1}{2^{l_i}}\right]^n = \sum_{i_1,\ldots,i_n} \frac{1}{2^{l_{i_1}}} \times \cdots \times \frac{1}{2^{l_{i_n}}}$$

The sum is over all combinations of values for $i_1, \ldots, i_n$ in $\{1, \ldots, I\}$.

Let's rewrite this in terms of possible values for $j = l_{i_1} + \cdots + l_{i_n}$:

$$K^n = \sum_{j=1}^{nl_I} \frac{N_{j,n}}{2^j}$$

$N_{j,n}$ is the number of sequences of $n$ codewords that have total length $j$. If the code is uniquely decodable, $N_{j,n} \leq 2^j$, so $K^n \leq nl_I$, which for big enough $n$ is possible only if $K \leq 1$.