

## Lecture 5: Variable Length Lossless Compression

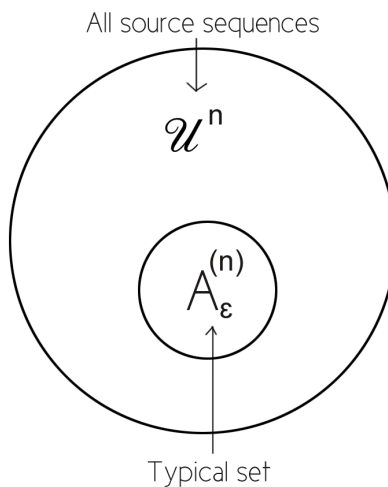
Lecturer: Tsachy Weissman

Scribe: Anthony Corso, Pengda Liu, Yuetong Wang

In this lecture we investigate strategies for variable length lossless compression. Several variable length encoding schemes are shown and their merits are discussed. Then a procedure for producing prefix codes for a dyadic source is presented. Lastly, Shannon Codes are presented as a generalization of the dyadic encoding procedure.

## 1 Review From Previous Class

$$U_1, U_2, \dots, U_n \sim \text{iid } U \in \mathcal{U} \quad (1)$$



**Figure 1:** Asymptotic Equipartition Property

$$P(U^n \in A_\epsilon^{(n)}) \approx 1 \quad (2)$$

$$\forall u^n \in A_\epsilon^{(n)} : p(u^n) \approx 2^{-nH(U)} \quad (3)$$

$$|A_\epsilon^{(n)}| \approx 2^{nH(U)} \quad (4)$$

This tells us that we need  $H(U)$  bits per source symbol for a near-lossless fixed length scheme. If fewer bits are used (say  $\alpha < H(U)$  bits), then the size of the set encoded by those  $\alpha$  bits is given by

$$|B_n| = 2^{n\alpha}$$

and we showed last time that the probability of a source sequence being within the set  $B_n$  goes to 0 and hence such a scheme makes an error with probability close to 1.

## 2 Variable Length Lossless Compression Examples

### 2.1 Example 1

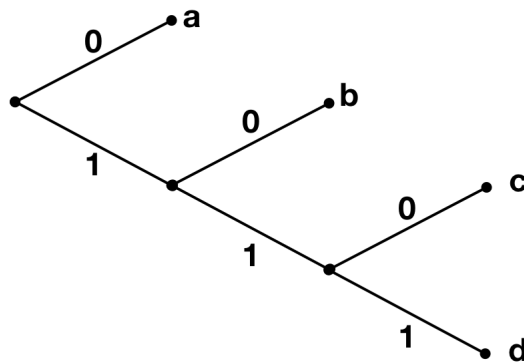
Given an alphabet with four letters in it

$$\mathcal{U} = \{a, b, c, d\}$$

The probability of each letter and a coding scheme is given in the table below:

$u$	$p(u)$	codeword $c(u)$	length $l(u)$
a	$1/2 = 2^{-1}$	0	1
b	$1/4 = 2^{-2}$	10	2
c	$1/8 = 2^{-3}$	110	3
d	$1/8 = 2^{-3}$	111	3

The code is  $\{c(u)\}_{u \in \mathcal{U}}$ . Note that any code can be represented by the nodes in a binary tree. For example, a binary tree representation of this code is shown in figure 2.



**Figure 2:** Binary tree representation

**Note:** For this code,  $l(u) = \log \frac{1}{p(u)}$ , therefore

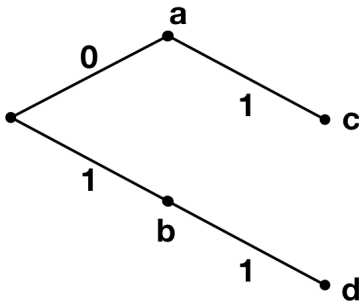
$$\bar{l} = \mathbb{E}[l(U)] = H(U) \quad (5)$$

Thus, this is a lossless scheme that achieves the ideal bits per source symbol, the entropy  $H(U)$  (we'll show in the next lecture that the entropy is fundamental limit of compression for variable length lossless schemes as well). Additionally, it is a prefix code which makes the encoding and decoding process very simple. Note that a code is a prefix code if and only if all the codewords are leaves in the binary tree representation of the code.

### 2.2 Example 2

Consider a new scheme for the same source (binary tree representation shown in figure 3).

$U$	codeword $c(U)$	length $l(u)$
a	0	1
b	1	1
c	01	2
d	11	3



**Figure 3:** The binary tree representation of the code

**Note:** The length of the code is  $\bar{l} < H(U)$ , i.e. better than the entropy. But, consider the encoding of three source symbols  $a b d$  and  $c b b$ :

$$a b d \xrightarrow{\text{encoded to}} 0111 \tag{6}$$

$$c b b \xrightarrow{\text{encoded to}} 0111 \tag{7}$$

Thus, different source sequences can have the same encoding, and this code cannot be decoded uniquely.

**Definition 1.** A code is **uniquely decodable (UD)** if every sequence of source symbols is mapped to a distinct binary representation.

**Definition 2.** A **prefix code** is a code where no codeword is the prefix of any other.

**Note:** A prefix code is uniquely decodable and can be decoded efficiently and on the fly (i.e. without needing entire binary sequence before decoding starts)

### 2.3 Exercise

Consider the code for the same source

$U$	codeword $c(U)$	length $l(u)$
a	10	2
b	00	2
c	11	2
d	110	3

Show that, though this code is not a prefix code, it is actually uniquely decodable. For the purpose of this exercise, you don't need to know the source distribution.

**Proof**

Part of HW 2.

□

### 3 Prefix code for dyadic distributions

**Definition 3.** A source is dyadic if

$$p(u) = 2^{-n_u} \quad (8)$$

where  $n_u$  is an integer  $\forall u \in \mathcal{U}$ .

Suppose we find a code such that

$$l(u) = n_u = \log \frac{1}{p(u)} \quad (9)$$

then

$$\bar{l} = \mathbb{E}[l(U)] = H(U) \quad (10)$$

Before proceeding to a prefix code for dyadic sources, we prove a technical lemma which says that the number of symbols with the lowest probability is even for a dyadic source. In particular, a dyadic source cannot have a single symbol with the lowest probability.

**Lemma 4.** Assume  $\mathcal{U}$  is dyadic with  $|\mathcal{U}| \geq 2$ , and let  $n_{\max} = \max_{u \in \mathcal{U}} n_u$ . Then the number of symbols  $u \in \mathcal{U}$  with  $n_u = n_{\max}$  is even.

**Proof:**

$$1 = \sum_u p(u) \quad (11)$$

$$= \sum_u 2^{-n_u} \quad (12)$$

$$= \sum_{n=1}^{n_{\max}} \left( \begin{array}{c} \# \text{ of symbols } u \\ \text{with } n_u = n \end{array} \right) 2^{-n} \quad (13)$$

$$(14)$$

Multiplying both sides by  $2^{n_{\max}}$ ,

$$\underbrace{2^{n_{\max}}}_{\text{even}} = \sum_{n=1}^{n_{\max}} \left( \begin{array}{c} \# \text{ of symbols } u \\ \text{with } n_u = n \end{array} \right) \cdot 2^{n_{\max}-n} \quad (15)$$

$$= \underbrace{\sum_{n=1}^{n_{\max}-1} \left( \begin{array}{c} \# \text{ of symbols } u \\ \text{with } n_u = n \end{array} \right) \cdot \underbrace{2^{n_{\max}-n}}_{\text{even}}}_{\text{even}} + \left( \begin{array}{c} \# \text{ of symbols } u \\ \text{with } n_u = n_{\max} \end{array} \right) \cdot 1 \quad (16)$$

Therefore, the number of letters  $u$  with  $n_u = n_{\max}$  must be even. ■

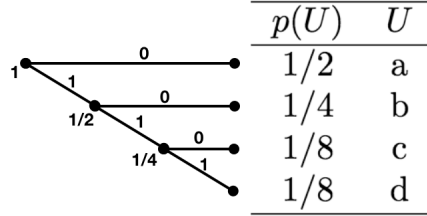
#### 3.1 Procedure for constructing prefix codes

Consider the following procedure:

- Choose 2 symbols with  $n_u = n_{\max}$  and merge them into a single symbol.
- We now have a new symbol with twice the probability.
- The new source is also dyadic.
- So repeat the first step until we are left with just one symbol.

**Note:** This procedure induces a binary tree and the codewords are the leaves of the constructed binary tree. As noted earlier, a code with all codewords on the leaf nodes is a prefix code.

### 3.2 Example of procedure



**Figure 4:** Procedure for generating a prefix code

**Note:** If  $p(u) = 2^{-n_u}$ , then the node  $u$  will participate in  $n_u$  merges (since each merge doubles the probability). Thus, the distance of  $u$  from the root is  $n_u$ , which is also equal to  $l(u)$ .

**Conclusion:** Our procedure yields a prefix code with

$$l(u) = n_u = \log \frac{1}{p(u)} \quad (17)$$

and, in particular,  $\bar{l} = H(U)$ .

## 4 Shannon Codes

For a general source let

$$n_u^* = \lceil \log \frac{1}{p(u)} \rceil \quad \forall u \in \mathcal{U} \quad (18)$$

**Note:**

$$\sum_{u \in \mathcal{U}} 2^{-n_u^*} = \sum_{u \in \mathcal{U}} 2^{-\lceil \log \frac{1}{p(u)} \rceil} \quad (19)$$

$$\leq \sum_{u \in \mathcal{U}} 2^{-\log \frac{1}{p(u)}} \quad (20)$$

$$= \sum_{u \in \mathcal{U}} p(u) = 1 \quad (21)$$

Consider a new source  $p^*(u) = 2^{-n_u^*}$ . While  $p^*(u)$  does not sum to 1 over  $\mathcal{U}$ , we can add new symbols to extend the source to  $\mathcal{U}^* \supseteq \mathcal{U}$  such that  $p^*(u)$  is dyadic over  $\mathcal{U}^*$  and  $\sum_{u \in \mathcal{U}^*} p^*(u) = 1$ .

Using the technique in the previous section, we can now construct a prefix code for the source  $p^*$  such that

$$l(u) = \log \frac{1}{p^*(u)} = n_u^* = \lceil \log \frac{1}{p(u)} \rceil \quad \forall u \in \mathcal{U} \quad (22)$$

Note that the code is defined over  $\mathcal{U}^*$  but we only care about the symbols in  $\mathcal{U}$ .

**Note:** This code is known as the *Shannon Code*.

The expected length for this code is

$$\bar{l} = \sum_u p(u)l(u) \tag{23}$$

$$= \sum_u p(u) \lceil \log \frac{1}{p(u)} \rceil \tag{24}$$

$$\leq \sum_u p(u) \left( \log \frac{1}{p(u)} + 1 \right) \tag{25}$$

$$= H(U) + 1 \tag{26}$$

So we are within 1 bit per source symbol of the entropy. To get even closer to the entropy, we can work with blocks of source symbols.

For the memoryless source  $U_1, U_2, \dots$  iid  $\sim \mathcal{U}$ , we can construct a Shannon code for  $U^N = (U_1, U_1, \dots, U_N)$ . Then

$$\frac{1}{N} \mathbb{E}[l(U^N)] \leq \frac{1}{N} (H(U^N) + 1) = H(U) + \frac{1}{N} \tag{27}$$

Thus as  $N$  becomes large, the bits used per source symbol tends toward the entropy.