

## Lecture 10: Feb 12

Lecturer: Akshay Krishnamurthy

Scribe: Dena Asta, Kirthevasan Kandasamy

**Disclaimer:** *These notes have not been subjected to the usual scrutiny reserved for formal publications. They may be distributed outside this class only with the permission of the Instructor.*

## 10.1 Codes

Codes are functions that convert strings over some alphabet into (typically shorter) strings over another alphabet. Recall, the coding problem

$$\mathcal{X} \ni x \longrightarrow \boxed{\text{Encoder}} \longrightarrow C(x) \in \Sigma^*$$

Here  $\Sigma$  is the dictionary (e.g. for binary codes  $\Sigma = \{0, 1\}$ ). Our goal is to have low Expected code length with respect to the distribution  $p$  of  $x$ ,  $\ell(C) \triangleq \mathbb{E}_{x \sim p} \ell(x)$  where  $\ell(x)$  is the length of  $C(x)$ .

### 10.1.1 Taxonomy of codes

Let  $X$  be a random variable taking values in a set  $\mathcal{X}$ . Let  $\Delta^*$  denote Kleene closure of the dictionary  $\Delta$  and  $C$  be the code – i.e  $C : \mathcal{X} \rightarrow \Sigma^*$ . The *extension* of the code  $C$  is a code of the form  $C : \mathcal{X}^* \rightarrow \Sigma^*$  defined by

$$C(x_1 x_2 \dots x_n) = C(x_1) C(x_2) \dots C(x_n), \quad n = 0, 1, \dots, x_1, x_2, \dots, x_n \in X.$$

Listed below are the types of codes we saw in class.

Codes	Description of $C$
Nonsingular	$C$ injective – i.e. $\forall x, x' \in \mathcal{X}, x \neq x' \implies C(x) \neq C(x')$
Uniquely Decodable	The extension of the code is nonsingular.
Prefix/Instantaneous/Self-punctuating	No code word prefixes another: for all distinct $x', x'' \in \mathcal{X}$ , $C(x')$ does not start with $C(x'')$

We illustrate this below with the following example for the symbols  $\{a, b, c, d\}$  taken from Chapter 5 in Cover and Thomas.

$X$	Singular	Nonsingular	Uniquely Decodable	Instantaneous
$a$	0	0	10	0
$b$	0	010	00	10
$c$	0	01	11	110
$d$	0	10	110	111

We begin with the following important results.

**Theorem 10.1 (Kraft-McMillan Inequality)** For any uniquely decodable code  $C : \mathcal{X} \rightarrow \Sigma$  where  $D = |\Sigma|$ ,

$$\sum_x D^{-\ell(x)} \leq 1. \quad (10.1)$$

Conversely, for all sets  $\{\ell(x)\}_{x \in \mathcal{X}}$  of numbers satisfying (12.1), there exists a prefix code  $C : \mathcal{X} \rightarrow \{1, 2, \dots, D\}^*$  such that  $\ell(x)$  is the length of  $C(x)$  for each  $x$ .

**Theorem 10.2**

1.  $H(X) \leq \ell(C)$  for all uniquely decodable codes.
2. For any  $\epsilon > 0$ , there exists large enough  $n$  and a code  $C_n : \mathcal{X}^n \rightarrow \Sigma^*$  such that  $\ell(C_n) \leq H(X) + \epsilon$ .

**Proof:** The proof for the first statemet follows by solving the convex program  $\min \sum_{x \in \mathcal{X}} p(x) \ell(x)$  subject to the constraints  $\sum D^{-\ell(x)} \leq 1$ .  $\sum_x p(x) = 1$ . For the second, we can use the Shannon code in blocks. That is, for an  $n$ -length sequence  $x_1^n = x_1 x_2 \dots x_n$  use the code lengths  $\ell_n(x_1^n) = \lceil -\log_D p(x_1^n) \rceil$  so that

$$H(X) = \frac{H(X_1^n)}{n} \leq \mathbb{E} \ell_n(X_1^n) \leq \frac{H(X_1^n) + 1}{n} = H(X) + \frac{1}{n}$$

■

**Proposition 10.3** The ideal codelengths for a prefix code with smallest expected codelength are

$$l^*(x) = \log_D \frac{1}{p(x)} \quad (\text{Shannon information content})$$

**Proof:** In last class, we showed that for all length functions  $l$  of prefix codes,  $E[l^*(x)] = H_p(X) \leq E[l(x)]$ .

■

While Shannon entropies are not integer-valued and hence cannot be the lengths of code words, the integers

$$\{\lceil \log_D \frac{1}{p(x)} \rceil\}_{x \in \mathcal{X}}$$

satisfy the Kraft-McMillan Inequality and hence there exists some uniquely decodable code  $C$  for which

$$H_p(x) \leq E[l(x)] < H_p(x) + 1, \quad x \in \mathcal{X} \quad (10.2)$$

by Theorem 12.1. Such a code is called **Shannon code**. Moreover, the lengths of code words for such a code  $C$  achieve the entropy for  $X$  asymptotically, i.e. if Shannon codes are constructed for strings of symbols  $x^n$  where  $n \rightarrow \infty$ , instead of individual symbols. Assuming  $X_1, X_2, \dots$  form an iid process, for all  $n = 0, 1, \dots$

$$\begin{aligned} H(X) &= \frac{H(X_1, X_2, \dots, X_n)}{n} \\ &\leq \frac{E[l(x_1, \dots, x_n)]}{n} \\ &< \frac{H(X_1, X_2, \dots, X_n)}{n} + \frac{1}{n} = H(X) + \frac{1}{n} \end{aligned}$$

by (12.2), and hence  $E[\frac{l(x_1, \dots, x_n)}{n}] \xrightarrow{n \rightarrow \infty} H(X)$ . If  $X_1, X_2, \dots$  form a startionary process, then a similar argument shows that  $E[\frac{l(x_1, \dots, x_n)}{n}] \xrightarrow{n \rightarrow \infty} H(\mathcal{X})$ , where  $H(\mathcal{X})$  is the entropy rate of the process.

**Theorem 10.4 (Shannon Source Coding Theorem)** *A collection of  $n$  iid random variables, each with entropy  $H(X)$ , can be compressed into  $nH(X)$  bits on average with negligible loss as  $n \rightarrow \infty$ . Conversely, no uniquely decodable code can compress them to less than  $nH(X)$  bits without loss of information.*

### 10.1.2 Non-singular vs. Uniquely decodable codes

Can we gain anything by giving up unique decodability and only requiring the code to be non-singular? First, the question is not really fair because we cannot decode sequence of symbols each encoded with a non-singular code easily. Second, (as we argue below) non-singular codes only provide a small improvement in expected codelength over entropy.

**Theorem 10.5** *The length of a non-singular code satisfies*

$$\sum_x D^{-l(x)} \leq l_{\max}$$

and for any probability distribution  $p$  on  $\mathcal{X}$ , the code has expected length

$$E[l(X)] = \sum_x p(x)l(x) \geq H_D(X) - \log_D l_{\max}.$$

**Proof:** Let  $a_l$  denote the number of unique codewords of length  $l$ . Then  $a_l \leq D^l$  since no codeword can be repeated due to non-singularity. Using this

$$\sum_x D^{-l(x)} = \sum_{l=1}^{l_{\max}} a_l D^{-l} \leq \sum_{l=1}^{l_{\max}} D^l D^{-l} = l_{\max}.$$

The expected codelength can be obtained by solving the following optimization problem:

$$\min \sum_x p(x)l(x) \quad \text{subject to} \quad \sum_x D^{-l_x} \leq l_{\max},$$

the convex non-singularity code constraint. Differentiating the Lagrangian  $\sum_x p_x l_x + \lambda \sum_x D^{-l_x}$  with respect to  $l_x$  and noting that at the global minimum  $(\lambda^*, l_x^*)$  it must be zero, we get :

$$p_x - \lambda^* D^{-l_x^*} \ln D = 0$$

which implies that  $D^{-l_x^*} = \frac{p_x}{\lambda^* \ln D}$ .

Using complementary slackness, noting that  $\lambda^* > 0$  for the above condition to make sense, we have :

$\sum_x D^{-l_x^*} = \sum_x \frac{p_x}{\lambda^* \ln D} = l_{\max}$  which implies  $\lambda^* = 1/(l_{\max} \ln D)$  and hence  $D^{-l_x^*} = p_x l_{\max}$ , or the optimum length  $l_x^* = -\log_D(p_x l_{\max})$ .

This gives the expected minimum codelength for nonsingular codes as  $\sum_x p_x l_x^* = -\sum_x p_x \log_D(p_x l_{\max}) = H_D(X) - \log_D l_{\max}$ . ■

In last lecture, we saw an example of a non-singular code for a process which has expected length below entropy. However, this is only true when encoding individual symbols. As a direct corollary of the above result, if symbol strings of length  $n$  are encoded using a non-singular code, then

$$E[l(X^n)] \geq H(X^n) - \log_D(n l_{\max})$$

Thus, the expected length per symbol can't be much smaller than the entropy (for iid processes) or entropy rate (for stationary processes) asymptotically even for non-singular codes, since the second term divided by  $n$  is negligible.

Thus, non-singular codes don't offer much improvement over uniquely decodable and prefix codes. In fact, the following result shows that any non-singular code can be converted into a prefix code while only increasing the codelength per symbol by an amount that is negligible asymptotically.

### 10.1.3 Huffman Coding

Is there a prefix code with expected length shorter than Shannon code? The answer is yes. The optimal (shortest expected length) prefix code for a given distribution can be constructed by a simple algorithm due to Huffman.

We introduce an optimal symbol code, called a *Huffman code*, that admits a simple algorithm for its implementation. We fix  $\Sigma = \{0, 1\}$  and hence consider binary codes, although the procedure described here readily adapts for more general  $\Sigma$ . Simply, we define the *Huffman code*  $C : X \rightarrow \{0, 1\}^*$  as the coding scheme that builds a binary tree from leaves up - takes the two symbols having the least probabilities, assigns them equal lengths, merges them, and then reiterates the entire process. Formally, we describe the code as follows. Let

$$\mathcal{X} = \{x_1, \dots, x_N\}, \quad p_1 = p(x_1), p_2 = p(x_2), \dots, p_N = p(x_N).$$

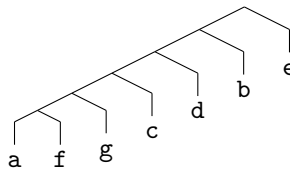
The procedure Huff is defined as follows:

```
Huff ( $p_1, \dots, p_N$ ):
  if  $N > 2$  then
     $C(1) \leftarrow 0, C(2) \leftarrow 1$ 
  else
    sort  $p_1 \geq p_2 \geq \dots p_N$ 
     $C' \leftarrow \text{Huff}(p_1, p_2, \dots, p_{N-2}, p_{N-1} + p_N)$ 
    for each  $i$ 
      if  $i \leq N - 2$  then  $C(i) \leftarrow C'(i)$ 
      else if  $i = N - 1$  then  $C(i) \leftarrow C'(N - 1) \cdot 0$ 
      else  $C(i) \leftarrow C'(N - 1) \cdot 1$ 
  return  $C$ 
```

For example, consider the following probability distribution:

symbol	a	b	c	d	e	f	g
$p_i$	0.01	0.24	0.05	0.20	0.47	0.01	0.02
Huffman code	000000	01	0001	001	1	000001	00001

The Huffman tree is build using the procedure described above. The two least probable symbols at the first iteration are 'a' and 'f', so they are merged into one new symbol 'af' with probability  $0.01 + 0.01 = 0.02$ . At the second iteration, the two least probable symbols are 'af' and 'g' which are then combined and so on. The resulting Huffman tree is shown below.



The Huffman code for a symbol  $x$  in the alphabet  $\{a, b, c, d, e, f, g\}$  can now be read starting from the root of the tree and traversing down the tree until  $x$  is reached; each leftwards movement suffixes a 0 bit and each rightwards movement adds a trailing 1, resulting in the code shown above in the table.

*Remark 1:* If more than two symbols have the same probability at any iteration, then the Huffman coding may not be unique (depending on the order in which they are merged). However, all Huffman codings on that alphabet are optimal in the sense they will yield the same expected codelength.

*Remark 2:* One might think of another alternate procedure to assign small codelengths by building a tree top-down instead, e.g. divide the symbols into two sets with almost equal probabilities and repeating. While intuitively appealing, this procedure is suboptimal and leads to a larger expected codelength than the Huffman encoding. You should try this on the symbol distribution described above.

*Remark 3:* For a  $D$ -ary encoding, the procedure is similar except  $D$  least probable symbols are merged at each step. Since the total number of symbols may not be enough to allow  $D$  variables to be merged at each step, we might need to add some dummy symbols with 0 probability before constructing the Huffman tree. How many dummy symbols need to be added? Since the first iteration merges  $D$  symbols and then each iteration combines  $D-1$  symbols with a merged symbols, if the procedure is to last for  $k$  (some integer number of) iterations, then the total number of source symbols needed is  $1 + k(D - 1)$ . So before beginning the Huffman procedure, we add enough dummy symbols so that the total number of symbols look like  $1 + k(D - 1)$  for the smallest possible value of  $k$ .

Now we will show that the Huffman procedure is indeed optimal, i.e. it yields the smallest expected codelength for any prefix code. Since there can be many optimal codes (e.g. flipping bits in a code still leads to a code with same codelength, also exchanging source symbols with same codelength still yields an optimal code) and Huffman coding only finds one of them, let's first characterize some properties of optimal codes.

Assume the source symbols  $x_1, \dots, x_N \in \mathcal{X}$  are ordered so that  $p_1 \geq p_2 \geq \dots \geq p_N$ . For brevity, we write  $l_i$  for  $l(x_i)$  for each  $i = 1, \dots, N$ . We first observe some properties of general optimal prefix codes.

**Lemma 10.6** *For any distribution, an optimal prefix code exists that satisfies:*

1. *if  $p_j > p_k$ , then  $l_j \leq l_k$ .*
2. *The two longest codewords have the same length and correspond to the two least likely symbols.*
3. *The two longest codewords only differ in the last two bits.*

**Proof:** The collection of prefix codes is well-ordered under expected lengths of code words. Hence there exists a (not necessarily unique) optimal prefix code. To see (1), suppose  $C$  is an optimal prefix code. Let  $C'$  be the code interchanging  $C(x_j)$  and  $C(x_k)$  for some  $j < k$  (so that  $p_j \geq p_k$ ). Then

$$\begin{aligned}
 0 &\leq L(C') - L(C) \\
 &= \sum_i p_i l'_i - \sum_i p_i l_i \\
 &= p_j l_k + p_k l_j - p_j l_j - p_k l_k \\
 &= (p_j - p_k)(l_k - l_j)
 \end{aligned}$$

and hence  $l_k - l_j \geq 0$ , or equivalently,  $l_j \leq l_k$ .

To see (2), note that if the two longest codewords had differing lengths, a bit can be removed from the end of the longest codeword while remaining a prefix code and hence have strictly lower expected length. An application of (1) yields (2) since it tells us that the longest codewords correspond to the least likely symbols. ■

We claim that Huffman codes are optimal, at least among all prefix codes. Because our proof involves multiple codes, we avoid ambiguity by writing  $L(C)$  for the expected length of a code word coded by  $C$ , for each  $C$ .

**Proposition 10.7** *Huffman codes are optimal prefix codes.*

**Proof:** Define a sequence  $\{\mathcal{A}_N\}_{N=2, \dots, |\mathcal{X}|}$  of sets of source symbols, and associated probabilities  $\mathcal{P}_N = \{p_1, p_2, \dots, p_{N-1}, p_N + p_{N+1} + \dots + p_{|\mathcal{X}|}\}$ . Let  $C_N$  denote a Huffman encoding on the set of source symbols  $\mathcal{A}_N$  with probabilities  $\mathcal{P}_N$ .

We induct on the size of the alphabets  $N$ .

1. For the base case  $N = 2$ , the Huffman code maps  $x_1$  and  $x_2$  to one bit each and is hence optimal.
2. Inductively assume that the Huffman code  $C_{N-1}$  is an optimal prefix code.
3. We will show that the Huffman code  $C_N$  is also an optimal prefix code.

Notice that the code  $C_{N-1}$  is formed by taking the common prefix of the two longest codewords (least-likely symbols) in  $\{x_1, \dots, x_N\}$  and allotting it to a symbol with expected length  $p_{N-1} + p_N$ . In other words, the Huffman tree for the merged alphabet is the merge of the Huffman tree for the original alphabet. This is true simply by the definition of the Huffman procedure. Let  $l_i$  denote the length of the codeword for symbol  $i$  in  $C_N$  and let  $l'_i$  denote the length of symbol  $i$  in  $C_{N-1}$ . Then

$$\begin{aligned} L(C_N) &= \sum_{i=1}^{N-2} p_i l_i + p_{N-1} l_{N-1} + p_N l_N \\ &= \underbrace{\sum_{i=1}^{N-2} p_i l'_i + (p_{N-1} + p_N) l'_{N-1}}_{L(C_{N-1})} + (p_{N-1} + p_N) \end{aligned}$$

the last line following from the Huffman construction. Suppose, to the contrary, that  $C_N$  were not optimal. Let  $\bar{C}_N$  be optimal (existence is guaranteed by previous Lemma). We can take  $\bar{C}_{N-1}$  to be obtained by merging the two least likely symbols which have same length by Lemma 12.6. But then

$$L(\bar{C}_N) = L(\bar{C}_{N-1}) + (p_{N-1} + p_N) \geq L(C_{N-1}) + (p_{N-1} + p_N) = L(C_N)$$

where the inequality holds since  $C_{N-1}$  is optimal. Hence,  $C_N$  had to be optimal. ■

**Remarks 10.8** *The numbers  $p_1, p_2, \dots, p_N$  need not be probabilities - just weights  $\{w_i\}$  taking arbitrary non-negative values. Huffman encoding in this case results in a code with minimum  $\sum_i p_i w_i$ .*

**Remarks 10.9** *Since Huffman codes are optimal prefix codes, they satisfy  $H(X) \leq E[l(X)] < H(X) + 1$ , same as Shannon code. However, expected length of Huffman codes is never longer than that of a Shannon code, even though for any given individual symbol either Shannon or Huffman code may assign a shorter codelength.*

**Remarks 10.10** *Huffman codes are often undesirable in practice because they cannot easily accomodate changing source distributions. We often desire codes that can incorporate refined information on the probability distributions of strings, not just the distributions of individual source symbols (e.g. English language.) The Huffman coding tree needs to be recomputed for different source distributions (e.g. English vs French).*

## 10.2 Connections to Machine Learning

Consider the classical statistical learning set up. We have data  $Z$  from some distribution  $p$ . We want to have good prediction on some learning task which is often characterised via,

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{Z \sim p} [L(Z, f(Z))]$$

where  $\ell$  is some loss function and  $R(f) = \mathbb{E}_{Z \sim p} [L(Z, f(Z))]$  is the risk. For example in linear regression  $Z = (X, Y)$  and  $L(Z, f(Z)) = (Y - f(X))^2$  where  $f(X) = \beta^\top X$  – then the objective becomes,  $\beta^* = \operatorname{argmin}_{\beta} \mathbb{E}_{X, Y \sim p} [(Y - \beta^\top X)^2]$ . However, the challenge is that we don't see the true distribution  $p$ . Instead we have samples  $Z_1^n = \{Z_1, \dots, Z_n\}$  where  $Z_i \sim p$ . A natural idea is to minimize the empirical risk,

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{Z \sim \hat{p}} [L(Z, f(Z))] = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n L(Z_i, f(Z_i))$$

Denote the empirical risk by  $R_n(f) = \mathbb{E}_{Z \sim \hat{p}} [L(Z, f(Z))]$  This procedure is called Empirical Risk Minimization (ERM). But does ERM work ?

We will begin this discussion by considering bounded loss functions  $\ell \in [0, 1]$ . We shall extend this in the next lecture. For this, we wish to bound the *excess risk*,  $R(\hat{f}) - R(f^*)$ . Noting that  $R_n(\hat{f}) \leq R_n(f^*)$ , we see that  $R(\hat{f}) - R(f^*) \leq |R(\hat{f}) - R_n(\hat{f})| + |R(f^*) - R_n(f^*)|$ . Therefore it is sufficient to have a uniform deviation bound of the form  $\forall f \in \mathcal{F}, |R_n(f) - R(f)| \leq \epsilon$ .

Using Hoeffding inequality we see that for any fixed  $f \in \mathcal{F}$ , with probability  $> 1 - \delta$

$$|R_n(f) - R(f)| \leq \sqrt{\frac{\log(2/\delta)}{2n}}$$

Therefore if  $\mathcal{F}$  is finite we see that with probability  $> 1 - \delta$

$$\forall f \in \mathcal{F} \quad |R_n(f) - R(f)| \leq \sqrt{\frac{\log(2|\mathcal{F}|/\delta)}{2n}} \quad (10.3)$$

which gives us the desired uniform deviation bound.

Note that we got the uniform deviation bound by assigning failure probability  $\delta/|\mathcal{F}|$  to each function and then using the union bound. But to apply the union bound it is sufficient if  $\mathcal{F}$  is countable. Suppose we have a distribution  $\zeta$  over  $\mathcal{F}$ . Then we can assign a failure probability of  $\delta(f) = \delta\zeta(f)$  for each  $f \in \mathcal{F}$  and apply the union bound. In particular, if we can assign a prefix code to the class then we can use the distribution  $\zeta(f) = 2^{-\ell(f)}$ .  $\sum_f \zeta(f) \leq 1$  by the Kraft inequality. This gives a deviation bound of the form,

$$|R_n(f) - R(f)| \leq \sqrt{\frac{(1 + \ell(f)) \log(2/\delta)}{2n}} \triangleq \epsilon(f)$$