

ECE 7680

Lecture AC – Arithmetic Coding

Objective: To introduce arithmetic coding

Introduction

Arithmetic coding overcomes some of the problems of Huffman coding, in particular the potential 1 bit surplus problem. It operates as a human might, using information already observed to predict what might be coming, and coding based on the prediction. In addition, the technique explicitly separates the prediction portion from the encoding portion.

In AC, a bit sequence is interpreted as an interval on the real line from 0 to 1. For example **01** is interpreted as $0.01\dots$, which corresponds (not knowing what the following digits are) to the interval $[0.01, 0.10)$ (in binary) which is $[0.25, 0.5)$ (base ten). (Make sure understanding on brackets.) A longer string **01101** corresponds to the interval $[0.01101, 0.01110)$. The longer the string, the shorter the interval represented on the real line.

Assume we are dealing with an alphabet $\mathcal{A} = \{a_1, \dots, a_I\}$, where a_I is a special symbol meaning “end of transmission.” The source produces the sequence $x_1, x_2, \dots, x_n, \dots$, and not necessarily i.i.d. We further assume (or model) that there is a predictor which computes, or estimates

$$P(x_n = a_i | x_1, x_2, \dots, x_{n-1}),$$

which is available at both encoder and decoder.

We divide the segment $[0, 1)$ into I intervals whose lengths are equal to the probabilities $P(x_1 = a_i), i = 1, 2, \dots, I$. The first interval is

$$[0, P(x_1 = a_1))$$

The second interval is

$$[P(x_1 = a_1), P(x_1 = a_1) + P(x_1 = a_2)),$$

and so forth. More generally, to provide for the possibility of considering other symbols than just x_1 , we define the lower and upper cumulative probabilities:

$$Q_n(a_i | x_1, \dots, x_{n-1}) = \sum_{j=1}^{i-1} P(x_n = a_j | x_1, \dots, x_{n-1})$$

$$R_n(a_i | x_1, \dots, x_{n-1}) = \sum_{j=1}^i P(x_n = a_j | x_1, \dots, x_{n-1})$$

Then, for example, a_2 corresponds to the interval $[Q_1(a_2), R_1(a_2))$.

Now we represent the probabilities for the next symbol. Take, for example, the interval for a_1 , and subdivide it into intervals $a_1 a_1, a_1 a_2, \dots, a_1 a_I$, so that the length of the interval for $a_1 a_j$ is *proportional* to $P(a_j | a_1)$. In fact, we take the length of the subinterval for $a_1 a_j$ to be

$$P(x_1 = a_1, x_2 = a_j) = P(x_1 = a_1)P(x_2 = a_j | x_1 = a_1)$$

Then we note that the sum of the lengths of these subintervals will be

$$\sum_j P(x_1 = a_1, x_2 = a_j) = P(x_1 = a_1),$$

which sure enough is the correct length.

More generally, we subdivide each of the intervals for $a_i a_j$ similarly to have length of

$$P(x_1 = a_i, x_2 = a_j) = P(x_1 = a_1)P(x_2 = a_j|a_1 = a_i).$$

Then, we continue subdividing each subinterval for strings of length N .

The following algorithm (Mackay, p. 151) shows how to compute the interval $[u, v)$ for the string $x_1 x_2, \dots, x_N$. (Note: this is for demonstration purposes, since it requires infinite precision arithmetic. In practice, the algorithm is arranged so that infinite precision is not required.)

Algorithm 1 Arithmetic coding

```

u = 0.0                                lower side of interval
v = 1.0                                upper side of interval
p = v - u                               length of current interval = probability of sequence
for n = 1 to N {
    Compute upper and lower cumulative probabilities
    Compute  $R_n(i|x_1, \dots, x_{n-1})$  for  $i = 1, 2, \dots, I$ 
    Compute  $Q_n(i|x_1, \dots, x_{n-1})$  for  $i = 1, 2, \dots, I$ 
     $u = u + pQ_n(x_n|x_1, \dots, x_{n-1})$            update lower side by joint length
     $v = u + pR_n(x_n|x_1, \dots, x_{n-1})$            update upper side by joint length
     $p = v - u$                                        compute new length = new probability
}

```

In encoding, the interval is subdivided for each new symbol. To encode the string x_1, x_2, \dots, x_N , we send the binary string whose interval lies *within* the interval determined by the sequence.

Example 1 (Mackay, p. 151) Suppose a bent coin with outcomes **a** and **b** is to be transmitted. We throw in the outcome \square to represent ‘end of transmission’.

Suppose we want to code **bbba** \square . We obtain the following table (the reasoning behind the probability values will be discussed below).

Context (sequence thus far)	Probability of next symbol		
	$P(\mathbf{a}) = 0.425$	$P(\mathbf{b}) = 0.425$	$P(\square) = 0.15$
b	$P(\mathbf{a} \mathbf{b}) = 0.28$	$P(\mathbf{b} \mathbf{b}) = 0.57$	$P(\square \mathbf{b}) = 0.15$
bb	$P(\mathbf{a} \mathbf{bb}) = 0.21$	$P(\mathbf{b} \mathbf{bb}) = 0.64$	$P(\square \mathbf{bb}) = 0.15$
bbb	$P(\mathbf{a} \mathbf{bbb}) = 0.17$	$P(\mathbf{b} \mathbf{bbb}) = 0.68$	$P(\square \mathbf{bbb}) = 0.15$
bbba	$P(\mathbf{a} \mathbf{bbba}) = 0.28$	$P(\mathbf{b} \mathbf{bbba}) = 0.57$	$P(\square \mathbf{bbba}) = 0.15$

The subdivision of the interval is portrayed in the following figure (Mackay, fig. 4.3)

When the first **b** is observed, the encoder knows that the encoded string will start as 01, 10, or 00 (see the figure), but doesn't know which. No output symbol is produced. For the next symbol **b**, the interval lies wholly within the interval 1, so the encoder outputs 1. However, no additional bits can be sent yet. At the third symbol **b**, there is still not quite enough information (but very close)! When the **a** is read, the interval lies entirely within 1001, which can now be output. When \square , we obtain the division shown on the right. The interval 10011101 is contained in the interval for **bbba** \square , so we send that sequence.

The **decoder** receives 100111101 and parses it one bit at a time. The probabilities are built up exactly as before. Once 10 have been parsed, it is known that the original string must have contained a **b**, since 10 lies within the interval **b**. Knowing this, the decoding computes $P(\mathbf{a}|\mathbf{b})$, $P(\mathbf{b}|\mathbf{b})$ and $P(\square|\mathbf{b})$, conditioned upon the symbol that has been received, then deduce the boundaries of the intervals **ba**, **bb** and **b□**. Eventually, the second **b** is decoded, which is used to condition the probabilities. This continues until decoding is complete. When the \square is reached the decoder knows that the end of file has been reached. \square

One of the benefits of arithmetic coding is that the worst case redundancy *for an entire bit string* (which may, for example, consist of an entire file) is **at most two bits**, assuming the probabilistic model is correct. Given a probabilistic model \mathcal{H} , the ideal message length for a sequence \mathbf{x} is $l(\mathbf{x}|\mathcal{H}) = -\log[P(\mathbf{x}|\mathcal{H})]$. Suppose that $P(\mathbf{x}|\mathcal{H})$ is *just barely* between two binary intervals. Then the next smaller binary intervals contained in $P(\mathbf{x}|\mathcal{H})$ are smaller by a factor of 4. This factor of 4 corresponds to $\log_2 4 = 2$ bits overhead worst case.

Probability models

Performance of the AC depends on having a good model for the source probabilities. The better the model, the better it might be expected that the code performs. In principle, any probabilistic model can be used. We mention here some useful concepts in developing one.

Suppose, as before, we deal with the case of independent events. We have outcomes **a**, **b**, and \square , with probabilities p_a, p_b and p_\square . Let l be the number of outcomes (number of coin tosses). p_a could be anywhere in the range $[0, 1]$, and we may not have any predisposition toward one value. We model this ambivalence by saying that

$$P(p_a) = 1 \quad \text{for } p_a \in [0, 1].$$

That is, it is uniformly distributed. This is a *prior probability*. If we had some predisposition about p_a , this could be incorporated into the prior model (using something like a β distribution, for example).

The whole point of Bayesian estimation (which is what we find we are talking about here) is to merge our prior inclinations in with the observations. This is a problem of inference, which we can state this way: given a sequence of F bits, of which F_a are **a**s and F_b are **b**s, infer p_a . The inference is accomplished by the posterior (“after”) — the probability of p_a *after* a measurement \mathbf{s} is made. We write

$$P(p_a|\mathbf{s}, F) = \frac{P(\mathbf{s}|p_a, F)P(p_a)}{P(\mathbf{s}|F)}.$$

Now why this? Well, we can write down the conditional probability in the numerator:

$$P(\mathbf{s}|p_a, F) = p_a^{F_a}(1 - p_a)^{F_b}$$

(describe why). As we have seen elsewhere, it seems that the conditioning is always easiest they way you don’t need it. We also find

$$P(\mathbf{s}|F) = \int P(\mathbf{s}|p_a, F)P(p_a)dp_a = \int p_a^{F_a}(1 - p_a)^{F_b}P(p_a)dp_a = \frac{\Gamma(F_a + 1)\Gamma(F_b + 1)}{\Gamma(F_a + F_b + 2)} = \frac{F_a!F_b!}{(F_a + F_b + 1)!}$$

So we could infer p_a as the most probable value (the maximizer) of the posterior. For example, we find $P(p_a|\mathbf{s} = \mathbf{aba}, F = 3) \propto p_a^2(1 - p_a)$, with maximum of $p_a = 2/3$. Or we could infer based on the mean, which is $3/5$.

We also want to be able to make predictions. Given a sequence \mathbf{s} of length F as evidence we find the prediction of drawing an \mathbf{a} as

$$P(\mathbf{a}|\mathbf{s}, F) = \int P(\mathbf{a}|p_a) P(p_a|\mathbf{s}, F) dp_a.$$

Note that in this case, we are using the entire posterior probability, so we incorporate all of our uncertainty about p_a . We also have $P(\mathbf{a}|p_a) = p_a$ (by its definition), so our predictor is

$$P(\mathbf{a}|\mathbf{s}, F) = \int p_a \frac{p_a^{F_a} (1 - p_a)^{F - F_a}}{P(\mathbf{s}|F)} dp_a = \frac{F_a + 1}{F_a + F_b + 2}.$$

This update rule is known as Laplace's rule, and is the rule that was used in the coder above.

We could write this as

$$P_L(a|x_1, \dots, x_{n-1}) = \frac{F_a + 1}{\sum_i (F_i + 1)}$$

Another model, known as the Dirichlet model, is more "responsive":

$$P_D(a|x_1, \dots, x_{n-1}) = \frac{F_a + \alpha}{\sum_i (F_i + \alpha)}$$

Typically, α is small, like 0.01.

This is not the only possible rule, and doesn't necessarily take into account the relationship that might exist between dependent variables.

Another application of AC

Mackay suggests another application as a way of generating a sample from a model, a sophisticated random number generator. As a simple-minded example, divide the unit interval into lengths p_i . Pick a point at random on the unit interval. The probability that your pin lies in interval i is p_i .

More generally, to generate a sample of bits from some random model, feed ordinary random bits into an arithmetic decoder for that model. This corresponds to picking a point at random in $[0, 1]$. The decoder will select a string at random from the modeled distribution, and will use the smallest number of random numbers to do the job.