

## Кодирование недвоичных данных при помощи бинаризации и двоичного арифметического кодирования

Допустим, надо закодировать с помощью двоичного арифметического кодера недвоичные данные.

Чтобы была возможность закодировать недвоичные данные с помощью двоичного кодера, необходимо произвести бинаризацию этих данных. Возникает закономерный вопрос: как бинаризовать, чтобы не было потери по сжатию?

Рассмотрим ансамбль  $X = \{A, B, C, D\}$ , с распределениями вероятностей  $p_A, p_B, p_C, p_D$ . Энтропия ансамбля:

$$H = -p_A \log p_A - p_B \log p_B - p_C \log p_C - p_D \log p_D.$$

Бинаризация равномерным кодом с независимым кодированием:

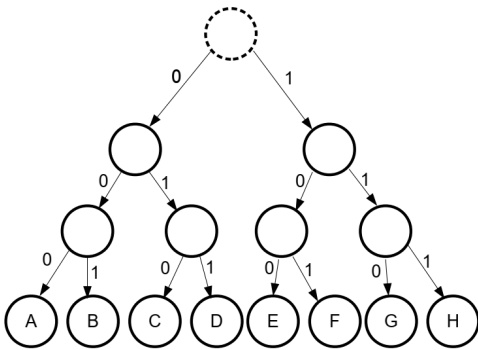
A → 00  
B → 01  
C → 10  
D → 11

т.е. возьмем эти коды и будем передавать в кодер по 1 биту  
в этом случае мы потеряем в сжатии, потому что результат сильно хуже начальной энтропии

При независимом кодировании первого и второго символов:

$$H_1 = H(X_1) + H(X_2) \geq H(X_1) + H(X_2|X_1) = H.$$

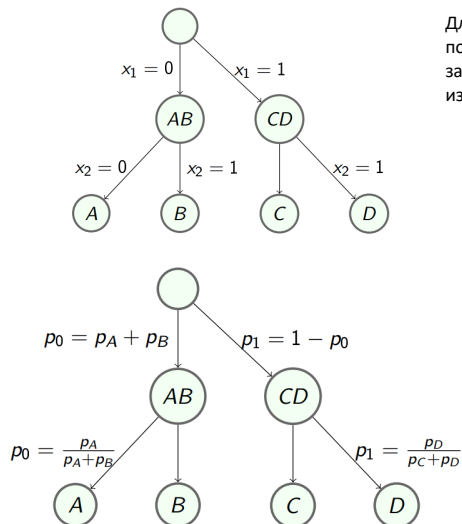
Древовидная бинаризация работает лучше, т.к. позволяет кодировать каждый бит с учетом предыдущего



Здесь, чтобы закодировать A, нужно передать в кодер три нуля по очереди. Но кодирование будет условное: то есть второй 0 будет кодироваться, учитывая условие, что первым тоже был 0.

Докажем, что при таком кодировании энтропия будет такая же, как и исходная, для примера выше

Рассмотрим ансамбль  $X = \{A, B, C, D\}$ , с распределениями вероятностей  $p_A, p_B, p_C, p_D$ .



Для буквы A сначала кодируем 0, показывая, что это AB, а не CD, и затем кодируем 0, показывая, что из A и B нам нужно A.

$$h_1 = -(p_A + p_B) \log(p_A + p_B) - (p_C + p_D) \log(p_C + p_D)$$

$$h_2 = (p_A + p_B) \left( -\frac{p_A}{p_A + p_B} \log \frac{p_A}{p_A + p_B} - \frac{p_B}{p_A + p_B} \log \frac{p_B}{p_A + p_B} \right) +$$

$$+ (p_C + p_D) \left( -\frac{p_C}{p_C + p_D} \log \frac{p_C}{p_C + p_D} - \frac{p_D}{p_C + p_D} \log \frac{p_D}{p_C + p_D} \right) =$$

$$H - h_1.$$

$$\Rightarrow h_1 + h_2 = H.$$

Вывод: при древовидной бинаризации мы не теряем в сжатии

$h_1$  - энтропия первого символа  
 $h_2$  - второго

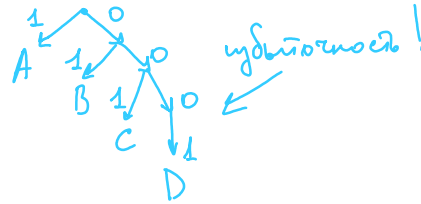
Унарная бинаризация

$n$	$unar(n)$
1	1
2	01
3	001

4	0001
5	00001
$n$	$\underbrace{00\dots 0}_{n-1}1$

Унарная бинаризация:

$A \rightarrow 1$   
 $B \rightarrow 01$   
 $C \rightarrow 001$   
 $D \rightarrow 0001$



$$\begin{aligned}
 h_1 &= -p_A \log p_A - (1 - p_A) \log(1 - p_A). \\
 h_2 &= -\frac{p_B}{p_B + p_C + p_D} \log \frac{p_B}{p_B + p_C + p_D} - \frac{p_C + p_D}{p_B + p_C + p_D} \log \frac{p_C + p_D}{p_B + p_C + p_D}. \\
 h_3 &= -\frac{p_C}{p_C + p_D} \log \frac{p_C}{p_C + p_D} - \frac{p_D}{p_C + p_D} \log \frac{p_D}{p_C + p_D}. \\
 h_4 &= 0.
 \end{aligned}$$

$$\begin{aligned}
 h_1 + (1 - p_A)h_2 + (1 - p_A - p_B)h_3 = \\
 -p_A \log p_A - (1 - p_A) \log(1 - p_A) - \\
 -p_B \log \frac{p_B}{p_B + p_C + p_D} - (p_C + p_D) \log \frac{p_C + p_D}{p_B + p_C + p_D} - \\
 -p_C \log \frac{p_C}{p_C + p_D} - p_D \log \frac{p_D}{p_C + p_D} = H.
 \end{aligned}$$

Что касается ренормализации, существует байтовая ренормализация для не двоичного алфавита, про которую препод в лекции толком ничего не сказал, поэтому я на экзамене тоже не скажу.

Здесь нет промежуточных вычислений, буфера и т.д. Здесь есть PUTBYTE, который сразу выдает байт на выход в файл

Байтовая ренормализация (range coder) для не двоичного<sup>3</sup> и двоичного<sup>4</sup> алфавита

```

1: while  $(L \oplus (L + R)) < 2^{24}$  or  $R < 2^{16}$  do
2:   if  $R < 2^{16} \wedge (L \oplus (L + R)) \geq 2^{24}$  then
3:      $R \leftarrow (!L + 1) \wedge (2^{16} - 1)$ 
4:   end if
5:   PUTBYTE  $(L \gg 24)$ 
6:    $R \leftarrow R \ll 8$ 
7:    $L \leftarrow L \ll 8$ 
8: end while

```

```

1: if  $(L \oplus (L + R)) < 2^{24}$  then
2:   PUTBYTE  $(L \gg 24)$ 
3:    $R \leftarrow R \ll 8$ 
4:    $L \leftarrow L \ll 8$ 
5: else if  $R < 2^{16}$  then
6:    $R \leftarrow (!L + 1) \wedge (2^{16} - 1)$ 
7:   PUTBYTE  $(L \gg 24)$ 
8:    $R \leftarrow R \ll 8$ 
9:    $L \leftarrow L \ll 8$ 
10: end if

```

главная байтовая ренормализация