

Современные проблемы информатики

Кодирование источника. Статистические коды

Фионов Андрей Николаевич

СибГУТИ

2020

Основные понятия

Код – последовательность символов некоторого (кодowego) алфавита, сопоставленная буквам или сообщениям источника

Двоичные коды строятся из символов алфавита $\{0, 1\}$

Кодовое слово – последовательность кодовых символов, сопоставленная некоторой букве алфавита источника

Длина кодowego слова – количество кодовых символов в слове

Кодовое пространство – множество, из которого выбираются кодовые слова

Виды кодов

Коды *фиксированной длины* – длины всех кодовых слов равны

Коды *переменной длины* – длины (хотя бы некоторых) кодовых слов различны

Разделимые (однозначно декодируемые) коды – коды, позволяющие восстановить (декодировать) символ источника единственным образом

Неразделимые коды – коды, не позволяющие декодировать символ однозначно

Префиксные коды – коды, в которых ни одно кодовое слово не является началом другого кодового слова

Полные коды – коды, в которых кодовые слова полностью покрывают всё кодовое пространство

Неполные коды – коды, оставляющие свободные (незанятые) места в кодовом пространстве

Примеры кодов

$$A = \{a, b, c, d, e\}$$

	1	2	3	4
<i>a</i>	000	0	00	1
<i>b</i>	001	1	01	10
<i>c</i>	010	10	100	10000
<i>d</i>	011	11	101	10001
<i>e</i>	100	100	11	10010

Коды фиксированной длины – 1 Коды переменной длины – 2, 3, 4

Разделимые коды – 1, 3, 4 Неразделимые коды – 2

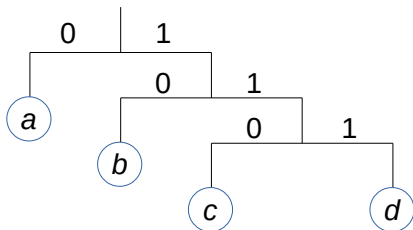
Префиксные коды – 1, 3

Полные коды – 2, 3 Неполные коды – 1, 4

Представление в виде дерева

$A = \{a, b, c, d\}$

a 0
 b 10
 c 110
 d 111

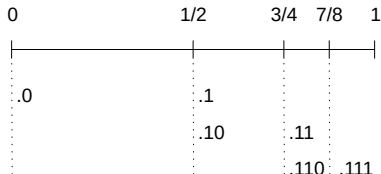


- полный код – всем листьям соответствуют буквы алфавита
- декодирование – проход по дереву

Интервальное представление

$$A = \{a, b, c, d\}$$

a	0	$[0, 1/2)$
b	10	$[1/2, 3/4)$
c	110	$[3/4, 7/8)$
d	111	$[7/8, 1)$



- разделимый код – интервалы не перекрываются
- декодирование – попадание в интервал

Неравенство Крафта–Макмиллана

1. Длины кодовых слов $\ell_1, \ell_2, \dots, \ell_N$ любого двоичного делимого кода удовлетворяют неравенству

$$\sum_{i=1}^N 2^{-\ell_i} \leq 1.$$

2. Для любого набора длин, удовлетворяющих этому неравенству, существует префиксный код с такими длинами кодовых слов.

Средняя длина кодового слова

$$\tilde{\ell} = p_1 \ell_1 + p_2 \ell_2 + \cdots + p_N \ell_N$$

Первая теорема Шеннона

Теорема о кодировании источника.

1. Средняя длина кодового слова разделимого кода не может быть меньше энтропии источника ($\tilde{\ell} \geq H$).
2. Наименьшая средняя длина кодового слова ($\tilde{\ell} = H$) достигается тогда и только тогда, когда каждый символ, появляющийся с вероятностью p , кодируется точно $-\log p$ битами.

Избыточность кода

$$r = \tilde{\ell} - H$$

Код Шеннона. Алфавитный код

Определение. *Кумулятивная вероятность* символа есть сумма вероятностей предшествующих символов

$$q_1 = 0, \quad q_i = \sum_{j=1}^{i-1} p_j, \quad i = 2, \dots, N$$

$$(q_1 = 0, \quad q_i = q_{i-1} + p_{i-1}, \quad i = 2, \dots, N)$$

Код Шеннона

Построение.

1. Упорядочиваем буквы алфавита по убыванию вероятностей.
2. Вычисляем кумулятивные вероятности.
3. В качестве кода буквы x берем $\lceil -\log p_x \rceil$ бит дробной части двоичного представления числа q_x .

Свойства кода Шеннона:

- Код является префиксным
- Избыточность $r < 1$ бит на символ

Код Шеннона

Декодируемость (префиксность).

Рассмотрим два кодовых слова для букв a_i и a_j с длинами $\ell_i, \ell_j, i < j$.

Вследствие сортировки $\ell_i \leq \ell_j$.

Из определения кумулятивных вероятностей $q_j - q_i \geq p_i$.

В соответствии с методом построения кода $\ell_i = \lceil -\log p_i \rceil \geq -\log p_i$.

Отсюда $p_i \geq 2^{-\ell_i}$ и $q_j - q_i \geq 2^{-\ell_i}$.

Значит i -е и j -е кодовые слова отличаются по меньшей мере в одном из ℓ_i разрядов. Код получается префиксным.

Избыточность.

$$\lceil -\log p_x \rceil < -\log p_x + 1$$

$$\begin{aligned}\tilde{\ell} &= \sum_x p_x(\lceil -\log p_x \rceil) \\ &< \sum_x p_x(-\log p_x) + \sum_x p_x \\ &= H(X) + 1\end{aligned}$$

Код Шеннона

Пример.

$$A = \{a, b, c, d, e\}, \mathbb{P} = (1/16, 3/16, 1/16, 4/16, 7/16), H = 1.97$$

$$1. A' = \{e, d, b, a, c\}, \mathbb{P} = (7/16, 4/16, 3/16, 1/16, 1/16)$$

$$2. \mathbb{Q} = (0, 7/16, 11/16, 14/16, 15/16) = (.0000, .0111, .1011, .1110, .1111)$$

$$3. e \rightarrow 00, d \rightarrow 01, b \rightarrow 101, a \rightarrow 1110, c \rightarrow 1111$$

$$\tilde{\ell} = 2.4375, r = 0.4675$$

Алфавитный код

Построение.

1. Вычисляем кумулятивные вероятности.
2. В качестве кода буквы x берем код наибольшего подынтервала, целиком входящего в интервал $[q_x, q_x + p_x)$.

Свойства алфавитного кода:

- Код является префиксным
- Избыточность $r < 2$ бит на символ
- Кодовые слова сохраняют порядок букв алфавита

Алфавитный код

Пример.

$$A = \{a, b, c, d, e\}, \mathbb{P} = (1/16, 3/16, 1/16, 4/16, 7/16), H = 1.97$$

$$1. \mathbb{Q} = (0, 1/16, 4/16, 5/16, 9/16) = (.0000, .0001, .0100, .0101, .1001)$$

a		b		c		d		e				
0	1/16			4/16	5/16			9/16				1
0	0	0	0	0	0	0	0	1	1	1	1	1
0	0	0	0	1	1	1	1	0	0	0	0	1
0	0	1	1	0	0	1	1	0	0	1	1	1
0	1	0	1	0	1	0	1	0	1	0	1	1

$$2. a \rightarrow 0000, b \rightarrow 001, c \rightarrow 0100, d \rightarrow 011, e \rightarrow 11$$

$$\tilde{\ell} = 2.6875, r = 0.7175$$

Код Хаффмана (Huffman)

Построение.

1. Находим два символа с наименьшими вероятностями, приписываем одному 0, другому 1.
2. Заменяем в алфавите найденные символы на один новый символ, суммируя их вероятности.

Указанные 2 шага повторяем до тех пор, пока в алфавите не останется один символ.

Код каждого исходного символа получаем, читая приписанные биты в обратном порядке.

Код Хаффмана

Свойства кода Хаффмана:

- Код является префиксным
- Код всегда полный
- Избыточность $r < 1$ бит на символ
- Код является *оптимальным* (самое главное свойство) – для заданного набора вероятностей невозможно построить другой код, имеющий меньшую среднюю длину кодового слова

Код Хаффмана

Пример.

$$A = \{a, b, c, d, e\}, \mathbb{P} = (1/16, 3/16, 1/16, 4/16, 7/16), H = 1.97$$

$$a : 0, c : 1, A_1 = \{ac, b, d, e\}, p(ac) = 2/16$$

$$ac : 0, b : 1, A_2 = \{acb, d, e\}, p(acb) = 5/16$$

$$acb : 0, d : 1, A_3 = \{acbd, e\}, p(acbd) = 9/16$$

$$acbd : 0, e : 1, A_4 = \{acbde\}, p(acbde) = 1 - \text{конец}$$

$$a \rightarrow 0000, b \rightarrow 001, c \rightarrow 0001, d \rightarrow 01, e \rightarrow 1$$

$$\tilde{\ell} = 2, r = 0.03$$

Уменьшение избыточности

Единственный способ – кодировать блоки символов.

Пусть m – размер блока. Если рассматривать блоки символов как буквы некоторого “супералфавита”, то можно достичь избыточности

$$r < 1/m \text{ бит на символ}$$

Пример.

$$A = \{a, b\}, p_a = 0.2, p_b = 0.8, m = 2$$

$$A^2 = \{aa, ab, ba, bb\}, p_{aa} = 0.04, p_{ab} = 0.16, p_{ba} = 0.16, p_{bb} = 0.64$$

Используем код Хаффмана для A^2

Избыточность на символ исходного алфавита A

$$r < 1/2$$

В общем случае

$$r < 1/m, \quad r \rightarrow 0 \text{ при } m \rightarrow \infty$$

Сложность блочного кодирования

Код Шеннона, код Хаффмана: нужно делать сортировку символов алфавита

Трудоемкость по времени T , объем памяти M при увеличении размера алфавита N :

$$T = O(N \log N), \quad M = O(N)$$

При кодировании блоков размера m

$$T = O(mN^m \log N), \quad M = O(N^m), \quad m \rightarrow \infty$$

Сложность блочного кодирования

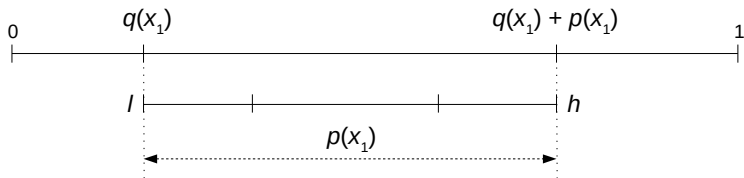
Алфавитный код: сортировку делать не нужно, требуется только вычислять кумулятивные вероятности

При кодировании блоков размера m

$$T = O(m^2), \quad M = O(m), \quad m \rightarrow \infty$$

Даёт начало наилучшему методу – арифметическому кодированию

Блочный алфавитный код



$$m = 4, X = x_1 x_2 x_3 x_4$$

$$1. \quad l = q(x_1), \quad h = l + p(x_1)$$

$$2. \quad l = q(x_1) + p(x_1)q(x_2), \quad h = l + p(x_1)p(x_2)$$

$$3. \quad l = q(x_1) + p(x_1)q(x_2) + p(x_1)p(x_2)q(x_3),$$

$$h = l + p(x_1)p(x_2)p(x_3)$$

$$4. \quad l = q(x_1) + p(x_1)q(x_2) + p(x_1)p(x_2)q(x_3) + p(x_1)p(x_2)p(x_3)q(x_4),$$

$$h = l + p(x_1)p(x_2)p(x_3)p(x_4)$$

Блочный алфавитный код

$$l = q(x_1) + p(x_1)q(x_2) + p(x_1)p(x_2)q(x_3) + p(x_1)p(x_2)p(x_3)q(x_4),$$
$$h = l + p(x_1)p(x_2)p(x_3)p(x_4)$$

t – число бит в представлении вероятностей

трудоемкость умножений: $2tt + 2(2t)t + 2(3t)t = 2t^2(1 + 2 + 3)$

для m символов $2t^2(1 + 2 + \dots + (m - 1)) = t^2m(m - 1) = O(m^2)$

Блочный алфавитный код

Пример.

$$p(a) = 1/16, p(b) = 3/16, p(c) = 1/16, p(d) = 4/16, p(e) = 7/16$$

$$q(a) = 0, q(b) = 1/16, q(c) = 4/16, q(d) = 5/16, q(e) = 9/16$$

$X = beed$

$$1. l = \frac{1}{16} \quad h = \frac{4}{16}$$

$$2. l = \frac{1}{16} + \frac{3}{16} \frac{9}{16} = \frac{43}{256} \quad h = \frac{43}{256} + \frac{3}{16} \frac{7}{16} = \frac{64}{256}$$

...

$$4. l = \frac{14767}{65536} \quad h = \frac{15355}{65536}$$

$$14767 = 0011 \ 1001 \ 1010 \ 1111$$

$$15354 = 0011 \ 1011 \ 1111 \ 1010$$

$$\text{код} = 00111010$$

Блочный алфавитный код – метод Рябко

$$l = q(x_1) + p(x_1)q(x_2) + p(x_1)p(x_2)q(x_3) + p(x_1)p(x_2)p(x_3)q(x_4),$$
$$h = l + p(x_1)p(x_2)p(x_3)p(x_4)$$

$$l = q(x_1) + p(x_1)q(x_2) + p(x_1)p(x_2) \cdot (q(x_3) + p(x_3)q(x_4)),$$
$$h = l + (p(x_1)p(x_2)) \cdot (p(x_3)p(x_4))$$

трудоемкость умножений: $4tt + 2(2t)(2t)$

$T = O(m \log^2 m \log \log m)$ (умножение методом Шёнхаге–Штрассена)

Арифметический код

L – левая граница интервала

R (range) – размер интервала

Вычисление границ интервала

$L \leftarrow 0, R \leftarrow 1;$

FOR ($i = 1, 2, \dots, m$) DO

$L \leftarrow L + Rq(x_i),$

$R \leftarrow Rp(x_i).$

Основные идеи:

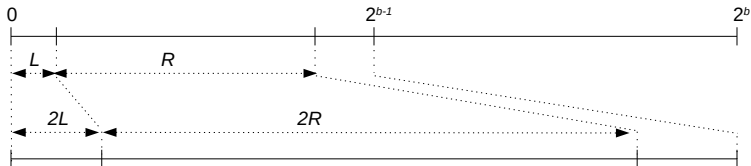
- 1 Переводим L, R в область целых чисел ($[0, 1) \rightarrow [0, 2^b)$)
- 2 При умножении на вероятность отсекаем дробную часть (немного жертвуем избыточностью)
- 3 Масштабируем интервал

Арифметический код – масштабирование

1. Интервал целиком лежит в левой половине

Выдаём на выход кодера бит 0

Масштабируем так:

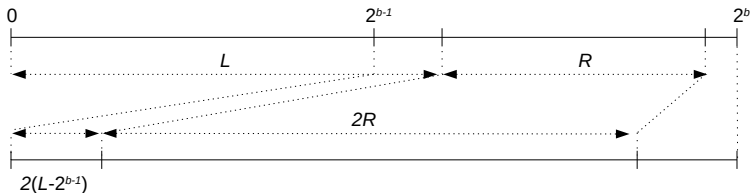


Арифметический код – масштабирование

2. Интервал целиком лежит в правой половине

Выдаём на выход кодера бит 1

Масштабируем так:

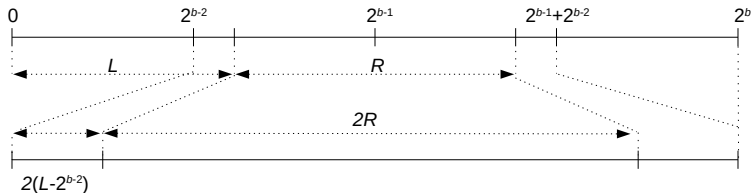


Арифметический код – масштабирование

3. Интервал целиком лежит в пределах 2 и 3 четверти

Увеличиваем на 1 счетчик S .

Масштабируем так:



В будущем при выдаче на выход кодера бита x после него надо будет выдать S противоположных бит

Арифметический код – масштабирование

В результате каждый бит выдаём на выход кодера с помощью следующего алгоритма

$\text{outS}(x)$

выдать бит x ,

S раз выдать бит $1 - x$,

$S \leftarrow 0$.

Арифметический код V1

Пусть p_1, p_2, \dots, p_N – целые, $\sum p_i = d$, т.е. $p(a_i) = p_i/d$.

Тогда $q_1 = 0, q_2, \dots, q_N, q_{N+1} = d$ тоже целые.

Полный интервал $[0, 2^b)$ – целочисленный.

Кодер: внутреннее состояние $L = 0, R = 2^{b-1}, S = 0$.

Процедура кодирования A_{enc} имеет три параметра – левая и правая границы интервала для символа, общий знаменатель вероятностей.

Для кодирования сообщения $x_1 x_2 \dots x_n$ вызываем

$A_{\text{enc}}(q(x_i), q(x_i) + p(x_i), d)$ для $i = 1, 2, \dots, n$. В ходе своего выполнения A_{enc} выдаёт биты на выход кодера и формирует новое внутренне состояние. В конце необходимо вывести код подынтервала внутри конечного состояния (максимум 2 бита).

Арифметический код V1

Aenc(l, h, d)

$L \leftarrow L + (R \times l) \operatorname{div} d$ (деление с отсечением дробной части)

$R \leftarrow (R \times (h - l)) \operatorname{div} d$

WHILE $R \leq 2^{b-2}$ DO (масштабирование)

IF $L + R \leq 2^{b-1}$ THEN outS(0)

ELSE IF $L \geq 2^{b-1}$ THEN outS(1), $L \leftarrow L - 2^{b-1}$

ELSE $S \leftarrow S + 1, L \leftarrow L - 2^{b-2}$

$L \leftarrow 2L, R \leftarrow 2R$

Завершение кодирования

$S \leftarrow S + 1$

IF $L \leq 2^{b-2}$ THEN outS(0)

ELSE outS(1)

Арифметический код V1

Декодер: внутреннее состояние V = первые b бит кода, $L = 0$, $R = 2^{b-1}$.

Декодирование символа

1. $t \leftarrow ((V - L + 1) \times d - 1) \text{ div } R$.
2. Находим символ s , для которого $l = q_s \leq t < q_{s+1} = h$.
- 3.

Adec(l, h, d)

$L \leftarrow L + (R \times l) \text{ div } d$

$R \leftarrow (R \times (h - l)) \text{ div } d$

WHILE $R \leq 2^{b-2}$ DO

IF $L + R \leq 2^{b-1}$ THEN ничего не делать

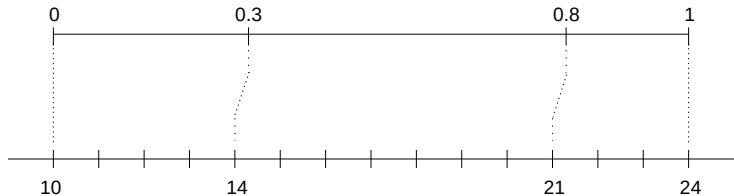
ELSE IF $L \geq 2^{b-1}$ THEN $L \leftarrow L - 2^{b-1}$, $V \leftarrow V - 2^{b-1}$

ELSE $L \leftarrow L - 2^{b-2}$, $V \leftarrow V - 2^{b-2}$

$L \leftarrow 2L$, $R \leftarrow 2R$, $V \leftarrow 2V$ + очередной бит кода

Арифметический код V1 – избыточность

Пример. $A = \{a, b, c\}$, $p_a = 0.3$, $p_b = 0.5$, $p_c = 0.2$



$$I(a) = -\log 0.3 = 1.74 \quad \ell_a = -\log(4/14) = 1.81$$

$$I(b) = -\log 0.5 = 1.00 \quad \ell_b = -\log(7/14) = 1.00$$

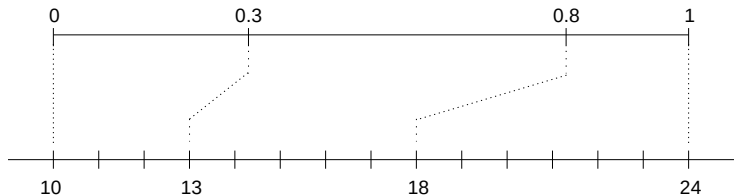
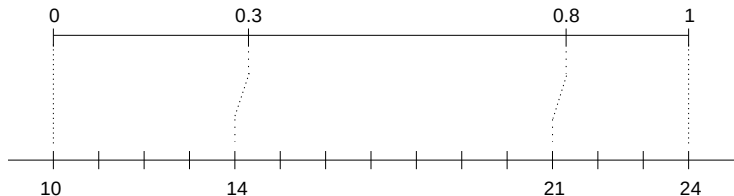
$$I(c) = -\log 0.2 = 2.32 \quad \ell_c = -\log(3/14) = 2.22$$

$$H = 1.4855 \quad \tilde{\ell} = 1.4867 \quad r = 0.0012$$

Арифметический код V2

Идея.

$$(R \times I) \operatorname{div} d \rightarrow (R \operatorname{div} d) \times I$$



Арифметический код V2

Кодер: внутреннее состояние $L = 0$, $R = 2^{b-1}$, $S = 0$.

Aenc2(l, h, d)

$k \leftarrow R \operatorname{div} d$

$L \leftarrow L + k \times l$

IF $h < d$ THEN $R \leftarrow k \times (h - l)$

ELSE $R \leftarrow R - k \times l$

масштабирование

Арифметический код V2

Декодер: внутреннее состояние D = первые b бит кода, $R = 2^{b-1}$
($D = V - L$).

Декодирование символа

1. $k \leftarrow R \operatorname{div} d$, $t \leftarrow \min(d - 1, D \operatorname{div} k)$.
2. Находим символ s , для которого $l = q_s \leq t < q_{s+1} = h$.
- 3.

Adec2(l, h, d)

$D \leftarrow D - k \times l$

IF $h < d$ THEN $R \leftarrow k \times (h - l)$

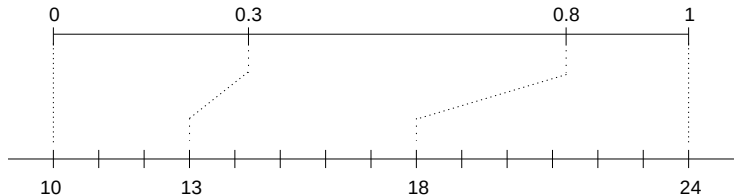
ELSE $R \leftarrow R - k \times l$

WHILE $R \leq 2^{b-2}$ DO

$R \leftarrow 2R$, $D \leftarrow 2D$ + очередной бит кода

Арифметический код V2 – избыточность

Пример. $A = \{a, b, c\}$, $p_a = 0.3$, $p_b = 0.5$, $p_c = 0.2$



$$I(a) = -\log 0.3 = 1.74 \quad \ell_a = -\log(3/14) = 2.22$$

$$I(b) = -\log 0.5 = 1.00 \quad \ell_b = -\log(5/14) = 1.49$$

$$I(c) = -\log 0.2 = 2.32 \quad \ell_c = -\log(6/14) = 1.22$$

$$H = 1.4855 \quad \tilde{\ell} = 1.6539 \quad r = 0.1684$$

Арифметический код – избыточность

$$r_{ac1} < N(\tau + \log e) \cdot 2^{-(b-2)} + 2/n$$

$$r_{ac2} < \log e \cdot 2^{-(b-2)+\tau} + 2/n$$

N – размер алфавита

τ – количество бит в представлении вероятности

b – битовый размер регистров кодера

n – длина сообщения

К О Н Е Ц

