

Блочное кодирование. Прямая и обратная теоремы кодирования

5 января 2023 г. 1:22

Код Гилберта-Мура лучше кода Шеннона отсутствием необходимости сортировать вероятности, но хуже с точки зрения избыточности. Применять такой код при побуквенном кодировании не имеет смысла как раз из-за большой избыточности. Чтобы уменьшить кодовую избыточность, можно использовать блочное кодирование. Чем больше длина блока, тем меньше избыточность на символ ($R = \bar{L} - H(X^n) = \bar{L} - nH_n(X)$)

- 1 Пусть $x \in X = \{0, 1, \dots, M-1\}$. Последовательность на выходе источника будем кодировать блоками $x = (x_1, x_2, \dots, x_n)$.
- 2 Каждый блок x длины n может рассматриваться как буква нового укрупнённого алфавита из всех комбинаций векторов длины n .
- 3 Применим любой алгоритм побуквенного кодирования для укрупнённого алфавита.

Энтропия на блок для укрупнённого алфавита вычисляется следующим образом:

Энтропия укрупнённого алфавита:

$$H(X^n) = - \sum_{x \in X^n} p(x) \log_2 p(x).$$

Далее применим побуквенное кодирование к блокам

Пусть r_n - избыточность укрупнённого алфавита, тогда $\bar{L} = H(X^n) + r_n$. Средние затраты на символ исходного алфавита:

$$\bar{R} = \frac{H(X^n) + r_n}{n} = \frac{H(X^n)}{n} + \frac{r_n}{n}.$$

избыточность
Харфман=1
Шеннон=1
Гилберт-Мур=2

← скорость кода, показывае, сколько бит мы тратим на символ

x - буква укрупнённого алфавита

вообще формула $\bar{L} \leq H+n$, а тут =, т.к. в общем виде

Если устремить n к бесконечности в формуле

то $\frac{r_n}{n}$ будет стремиться к уменьшению

а $\frac{H(X^n)}{n}$ - ни что иное, как энтропия на символ (в данном случае на блок, т.к. каждый символ нового алфавита - по сути блок)

- Для источника без памяти $H(X^n) = nH(X)$ получим:

$$\bar{R} = H(X) + \frac{r_n}{n}.$$

Если $n \rightarrow \infty$, то $\frac{r_n}{n} \rightarrow 0$.

- Для источника с памятью $H(X^n) \leq nH(X)$, и поэтому $\frac{H(X^n)}{n} \leq H(X)$

$$\lim_{n \rightarrow \infty} \frac{H(X^n)}{n} = H_\infty(X)$$

Прямая и обратная теоремы кодирования

Обратная теорема кодирования

$\bar{R} \geq H_\infty(X)$ (средняя длина на символ \geq энтропии, которая вычисляется при устремлении размера блока к бесконечности)

Прямая теорема кодирования

Можно найти такой префиксный код, для которого будет верно неравенство $\bar{R} \leq H_\infty(X) + \epsilon$, где ϵ - какое-то небольшое положительное число. То есть мы устремляем n к ∞ и можем достигнуть этого ϵ .

Доказательство

1. Берем исходный алфавит и укрупняем его. Получаем новый алфавит, в котором каждая буква - блок длиной n символов. Понятно, что средняя длина кодового слова будет больше или равна энтропии блока из n символов. Т.е.: $x \in X^n$, $\bar{L}_n \geq H(X^n)$
2. Для любого n существует код (например, Шеннона), такой что $\bar{L}_n \leq H(X^n) + 1$
3. Вычисляем среднюю длину на символ, получаем, что она больше или равна энтропии при устремлении n к ∞ : $R_n = \frac{\bar{L}_n}{n} \geq \frac{H(X^n)}{n} \geq H_\infty(X)$ для всех n (обратная теорема доказана)
4. Если мы в обратную сторону посмотрим, то можно найти такое большое n , что избыточность, которую мы прибавляем, будет стремиться к значению $H_\infty(X) + \epsilon$. Т.е.: для $n \geq n_0$ и $\epsilon > 0$ $R_n \leq H_\infty(X) + \epsilon$

Этот подход позволяет сжать файл близко к энтропии

- 1 Код Харфмана для укрупнённого алфавита сложно использовать на практике. Если $|X| = 256$, то для $n = 2$ $|X^n| = 65536$.
- 2 Код Шеннона сложно использовать, так как он требует сортировки.
- 3 Код Гилберта-Мура хорошо подходит для кодирования блоков.

т.к. для кода Шеннона $\bar{L} \leq H+1$, где 1 - это ϵ , т.е. избыточность

\geq , т.к. $\bar{L}_n \geq H(X^n)$

для $n=2$ реально реализовать, а для $n=10$, например, памяти не хватит (новый алфавит не влезет в переменную?)

Большинство алгоритмов проблематично сертифицировать реализуется, но большая избыточность

- Арифметическое кодирование является обобщением кода Гилберта-Мура для случая блочного кодирования.