

Оптимальное кодирование.

Лекция 7.
9.06

$\exists \varphi: A \rightarrow B^*$ — алфавитное кодирование
с длинами код. слов $len(\varphi(a_i)) = l_i, 1 \leq i \leq m,$
 φ^* кодирует ДИС без памяти $(A, \vec{p}), \vec{p} = (p_1, \dots, p_m)$
Хотим: $len(\varphi^*(A^n)) \rightarrow \min$ при $n \geq 1$

т.к. все A_i одинаково распределены,

$$M(len(\varphi^*(A^n))) = n \cdot M(len(\varphi^*(A_1))) = n \cdot \sum_{i=1}^m p_i l_i$$

ОПР. Пусть $\varphi: A \rightarrow B^*$ — алф. кодир. из $A = \{a_1, \dots, a_m\}$ в $B = \{b_1, \dots, b_r\}$

$$\vec{p} = (p_1, \dots, p_m)$$

ДИС

Средней длиной кода $\varphi(A)$
назов. $l^\varphi := \sum_{i=1}^m p_i l_i$

ОПР. Алфавитное кодир $\varphi: A \rightarrow B^*$ и код $\varphi(A)$ наз. оптимальным,
если φ — однозначно декодируемо и
средняя длина l^φ минимальна.

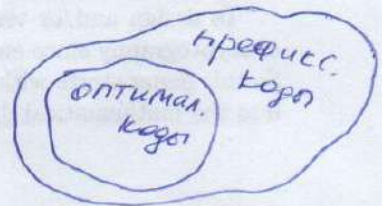
Замеч. 1) Если оптимальное кодир существует \Rightarrow

\Rightarrow существ. префиксное кодир с таким же набором код. слов
и это префикс. кодир — тоже оптимальное

2) т.к. $l_i \geq 1 \forall i \Rightarrow l^\varphi \geq 1$

так, при $m \leq \infty$ оптимальной код трив:

$$\varphi(a_i) = b_i; 1 \leq i \leq m$$



УТВ. Оптимальное кодир φ — существует.

факт: ДИС $(A, \vec{p}) \Rightarrow l^\varphi = \sum_{i=1}^m p_i l_i$ — функция от арг. (A, \vec{p})

$\Rightarrow \exists \inf_{\varphi} l^\varphi$ этой функции.

Докажем, что $\exists \varphi: A \rightarrow B^*$, т.ч. $l^\varphi = \inf_{\varphi} l^\varphi$

Л-во. $\exists L \in \mathbb{N}$ — наименьшее, т.ч. $m \leq \infty$. Тогда $\exists \varphi$ -равномер: $len \varphi(a_i) = L \forall 1 \leq i \leq m$

$$\Rightarrow \inf_{\varphi} l^\varphi \leq L$$

Не будем рассматрив. кодир-я φ , для кот. $l^\varphi > L$.

$$l^\varphi = \sum_{i=1}^m p_i l_i \leq L$$

Если $p_i = 0 \Rightarrow l_i$ — не вносит на l^φ . Удаляем эти слагаемые.

Если $p_i > 0 \Rightarrow l_i \leq \frac{L}{p_i} \leq \frac{L}{p}$, где $p = \min\{p_1, \dots, p_m\}$.

Нам нужно найти кодирование, у кот. все $l_i \leq \frac{L}{p}$ (при $p_i > 0$)

а l_i при $p_i = 0$ не вноят

l^φ при таких ограничениях
может прин. лишь конеч. число различных значений

\Rightarrow при некотором φ достигает \min .

III) Если алф. кодир. φ - однозначно декодируемо, то

$$L^\varphi \geq \frac{H(\vec{p})}{\log_2 \Phi}, \text{ причем равенство достигается тогда и т. тогда,}$$

когда все $\begin{cases} p_i = \Phi^{-l_i} > 0 \\ p_i = 0 \end{cases}$

Ф.во. $H(\vec{p}) - L^\varphi \log_2 \Phi =$

$$= - \sum_{i=1}^m p_i \log_2 p_i - \log_2 \Phi \sum_{i=1}^m p_i l_i = \left. \begin{matrix} \text{раскр. скобки,} \\ l_i \text{- лог. логарифм} \end{matrix} \right\}$$

$$= - \sum_{i=1}^m p_i \log_2 p_i + \sum_{i=1}^m p_i \log_2 \Phi^{-l_i} = \sum_{i=1}^m p_i \log_2 \frac{\Phi^{-l_i}}{p_i} \leq \left. \begin{matrix} \log_2 p_i \leq \log_2 (p_i - 1) \\ \text{Вентри,} \\ \text{McMillan} \end{matrix} \right\} \leq$$

$$\leq \log_2 e \cdot \sum_{i=1}^m p_i \left(\frac{\Phi^{-l_i}}{p_i} - 1 \right) = \log_2 e \left(\sum_{i=1}^m \Phi^{-l_i} - \sum_{i=1}^m p_i \right) \leq 0$$

III) Если распред. \vec{p} - невырожденное, то существует такое префикс. кодир. φ , для кот. справедливо:

$$L^\varphi < 1 + \frac{H(\vec{p})}{\log_2 \Phi}$$

Следствие. Среднее длина оптимального алф. кодирования φ удовл. нерав.-вам:

$$\frac{H(\vec{p})}{\log_2 \Phi} \leq L^\varphi < 1 + \frac{H(\vec{p})}{\log_2 \Phi}$$

(Без г-ва)

Осн. алгоритмы для построения префиксных кодов:

алг. Фано и Хатфманя

строит код малой средней длины.

строит оптимальный код

(не без. оптимальный)

Алгоритм Фано ($\Phi=2$)

Вход: ДЦС $\mathcal{A} = \{a_1, \dots, a_m\}$, $\vec{p} = (p_1, \dots, p_m)$

- алф. упорядочен так, чтобы $p_1 \geq p_2 \geq \dots \geq p_m$ (!)

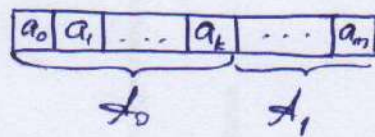
Кодовой алфавит: $\mathcal{B} = \{0, 1\}$ $\Phi = \# \mathcal{B} = 2$.

Выход: кодирование $\varphi: \mathcal{A} \rightarrow \mathcal{B}^*$.

Алгоритм: Шаг 1. Выберем $k \in [1, m]$:

$$\left| \sum_{i=1}^k p_i - \sum_{i=k+1}^m p_i \right| \longrightarrow \min$$

Шаг 2. Разбиваем A на подмнож.



$$A = A_0 \cup A_1$$

$$A_0 := \{a_1, \dots, a_k\} \quad A_1 := \{a_{k+1}, \dots, a_m\}$$

Шаг 3. Имеем множ. $A_{i_1, \dots, i_s} \subseteq A$; $i_1, \dots, i_s \in B$

Если $A_{i_1, \dots, i_s} = \{a_j\} \rightarrow$ присвоим симв. a_j код $i_1 \dots i_s$:
 $\varphi(a_j) := i_1 \dots i_s$.

Иначе: $A_{i_1, \dots, i_s} = \{a_j, \dots, a_l\}$

Выберем $k \in [j, l]$:

$$\left| \sum_{i=j}^k p_i - \sum_{i=k+1}^l p_i \right| \longrightarrow \min$$

Шаг s+1. Разбиваем $A_{i_1, \dots, i_s} := A_{i_1, \dots, i_s, 0} \cup A_{i_1, \dots, i_s, 1}$.

Пример. $A = \{a, b, c, d, e\}$

$$\vec{p} = \{0,3; 0,2; 0,2; 0,2; 0,1\}$$

$$\phi = 2$$

$$B = \{0, 1\}$$

(a) 0,3	0	0		$\varphi(a) = 00$
(b) 0,2	0	1		$\varphi(b) = 01$
(c) 0,2	1	0		$\varphi(c) = 10$
(d) 0,2	1	1	0	$\varphi(d) = 110$
(e) 0,1	1	1	1	$\varphi(e) = 111$

Средняя длина кода:

$$L^{\varphi} = (0,3 + 0,2 + 0,2) \cdot 2 + (0,2 + 0,1) \cdot 3 = 2,3 \text{ (бит)}.$$

Алгоритм Хаффмана ($\Phi=2$)

Вход: ДЧС $A = \{a_1, \dots, a_m\}$ $\vec{p} = (p_1, \dots, p_m)$
 - алф. упорядочен так, чтобы $p_1 \geq p_2 \geq \dots \geq p_m$ (!)
 Кодовый алфавит: $B = \{0, 1\}$ $\Phi = \#B = 2$.

Выход: Кодирование $f: A \rightarrow B^*$.

Алгоритм: ① Построим корневое разлн. дерево $G = (A^{(i)}, U)$
 $A^{(i)}$ - множ. вершин, U - множ. ребер.

Положим $A^{(0)} := A$ - листьями в дереве G .
 (самый низкий уровень).

$A^{(i)}$ - вершины предпоследнего уровня. Для их построения выполняем редукцию $A^{(i)}$:

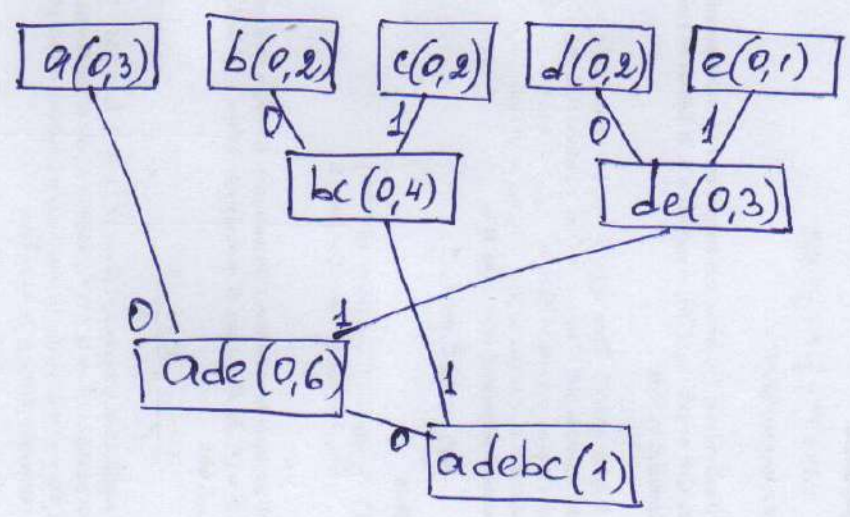
$A^{(1)} := \{a_1, \dots, a_{m-2}, a_{m-1} \cup a_m\}$ - две вершины с мин. вероятностями заменено на одну.
 $\vec{p}^{(1)} = (p_1, \dots, p_{m-2}, p_{m-1} + p_m)$

~~Вывод~~ Выполняем редукцию множ. $A^{(s)}$ и строим $A^{(s+1)}$, пока не получим $A^{(k)} = \{a_1 \cup \dots \cup a_m\}$
 $p^{(k)} = (1)$, k - высота дерева

② Код. слово $f(a_i)$ - запись последовательно метки ребер на пути от корня $a_1 \cup \dots \cup a_m$ к листу a_i .

Пример.

- $f(a) = 00$
- $f(b) = 10$
- $f(c) = 11$
- $f(d) = 010$
- $f(e) = 011$



Сред. длина $l^f = 2,3$ бит.