

Дерево оптимального поиска

Пусть двоичное дерево поиска содержит ключи $k_1, k_2, k_2, \dots, k_n$.

Задача № 1. Пусть каждый пользователь разыскивает только один из имеющихся в дереве ключей, причем, известна вероятность того, что пользователя интересует именно i -й ключ:

$$p_i = P|_{x=k_i}.$$

То, что пользователя интересуют только имеющиеся ключи, выражается равенством

$$\sum_{i=1}^n p_i = 1.$$

Требуется построить такое дерево, чтобы

$$P_I = \sum_{i=1}^n p_i h_i \rightarrow \min ,$$

h_i – уровень вершины i (число, на единицу большее глубины).

Задача № 2. Пусть каждый пользователь разыскивает ключ, который может либо присутствовать, либо отсутствовать в дереве. Вероятность успешного поиска сформулирована выше, а вероятность неудачи есть

$$q_i = P|_{k_i < x < k_{i+1}}, \quad i = 1, 2, 3, \dots, n-1,$$

$$q_0 = P|_{x < k_1}, \quad q_n = P|_{x > k_n}$$

То, что в дереве могут быть только интересующие или не интересующие пользователя ключи, выражается равенством

$$\sum_{i=1}^n p_i + \sum_{j=0}^n q_j = 1.$$

Требуется построить такое дерево, чтобы

$$P_{\Sigma} = P_I + P_E = \sum_{i=1}^n p_i h_i + \sum_{j=0}^n q_j h_j' \rightarrow \min$$

p_i – вероятность обращения к i -му ключу в дереве;

q_i – вероятность обращения к отсутствующим ключам, «расположенным» между i -м и $(i-1)$ -м ключами в дереве;

q_0 – вероятность обращения к отсутствующим ключам, «меньшим», чем 0-й (по номеру) ключ в дереве;

q_n – вероятность обращения к отсутствующим ключам, «большим», чем n -й (по номеру) ключ в дереве;

h_j' – уровень специальной вершины j .

Разумеется, каждая из этих задач может быть решена. Однако, величины p_i и q_j не могут быть известны точно. Они могут быть оценены приближенно, на основе статистики обращения пользователей к дереву поиска:

$$p_m = \frac{a_m}{\sum_{i=1}^n a_i + \sum_{j=0}^n b_j}, \quad q_m = \frac{b_m}{\sum_{i=1}^n a_i + \sum_{j=0}^n b_j},$$

где:

a_i – количество обращений к i -му ключу в дереве,

b_i – количество обращений к несуществующим ключам, «расположенным» между i -м и $(i-1)$ -м ключами в дереве,

b_0 – количество обращений к несуществующим ключам, «меньшим», чем 0-й (по номеру) ключ в дереве,

b_n – количество обращений к отсутствующим ключам, «большим», чем n -й ключ в дереве.

Задача № 3. Требуется построить такое дерево, чтобы

$$P = P_I' + P_E' = \sum_{i=1}^n a_i h_i + \sum_{j=0}^n b_j h_j' \rightarrow \min$$

Функция, которая минимизируется в задаче № 3, пропорциональна функции в задаче № 2. Коэффициент пропорциональности

$$\sum_{i=1}^n a_i + \sum_{j=0}^n b_j$$

от вида дерева не зависит.

Утверждение.

$$P = P_L + W + P_R, \quad \text{где} \quad W = \sum_{i=1}^n a_i + \sum_{j=0}^n b_j.$$

P называется *длиной пути* дерева, W называется *весом* дерева.

Доказательство.

$$\begin{aligned} & \sum_{i=1}^n a_i h_i + \sum_{j=0}^n b_j h_j' = \\ &= \sum_{i=1}^{m-1} a_i h_i + a_m \underbrace{h_m}_{=1} + \sum_{i=m+1}^n a_i h_i + \sum_{j=0}^{m-1} b_j h_j' + \sum_{j=m}^n b_j h_j' = \\ &= \sum_{i=1}^{m-1} a_i h_i + \sum_{j=0}^{m-1} b_j h_j' + \sum_{i=m+1}^n a_i h_i + \sum_{j=m}^n b_j h_j' + a_m = \\ & \left(\underbrace{\sum_{i=1}^{m-1} a_i (h_i - 1) + \sum_{j=0}^{m-1} b_j (h_j' - 1)}_{=P_L} \right) + \sum_{i=1}^{m-1} a_i + \sum_{j=0}^{m-1} b_j + \\ &+ \left(\underbrace{\sum_{i=m+1}^n a_i (h_i - 1) + \sum_{j=m}^n b_j (h_j' - 1)}_{=P_R} \right) + \sum_{i=m+1}^n a_i + \sum_{j=m}^n b_j + a_m = \\ &= P_L + P_R + \underbrace{\sum_{i=1}^n a_i + \sum_{j=0}^n b_j}_{=W} \end{aligned}$$

Замечание.

$$P_{\Sigma} = \frac{P}{W}.$$

Алгоритм построения дерева

Пусть T_{ij} – оптимальное поддереву, составленное из ключей $k_{i+1}, k_{i+2}, \dots, k_j$ (ключ k_i в дереве отсутствует).

Пусть вес этого дерева – w_{ij} , целевая функция (длина пути) – p_{ij} .

Ясно, что $w_{0n} = W$, $p_{0n} = P$.

$$w_{ij} = \sum_{s=i+1}^j a_s + \sum_{t=i}^j b_t, \quad p_{ij} = \sum_{s=i+1}^j a_s h_s + \sum_{t=i}^j b_t h'_t.$$

Тогда:

$$w_{ii} = b_i, \quad 0 \leq i \leq n,$$

$$w_{ij} = w_{i,j-1} + a_j + b_j, \quad 0 \leq i < j \leq n,$$

$$p_{ii} = w_{ii}, \quad 0 \leq i \leq n,$$

$$p_{ij} = w_{ij} + \min_{k: i < k \leq j} (p_{i,k-1} + p_{kj}), \quad 0 \leq i < j \leq n.$$

r_{ij} – то значение k , при котором достигается этот **min**.

Пример

Пусть $n = 4$, а числа пользовательских обращений представлены в таблице:

i	0	1	2	3	4
a_i		20	10	5	3
b_i	64	32	16	8	4

Без затруднений и объяснений заполняется таблица

$w_{00} = 64$	$w_{01} = 64 + 52 = 116$	$w_{02} = 116 + 26 = 142$	$w_{03} = 142 + 13 = 155$	$w_{04} = 155 + 7 = 162$
	$w_{11} = 32$	$w_{12} = 32 + 26 = 58$	$w_{13} = 58 + 13 = 71$	$w_{14} = 71 + 7 = 78$
		$w_{22} = 16$	$w_{23} = 16 + 13 = 29$	$w_{24} = 29 + 7 = 36$
			$w_{33} = 8$	$w_{34} = 8 + 7 = 15$
				$w_{44} = 4$

Формирование величин p_{ij} , r_{ij} несколько сложнее и требует пояснений.

Разница индексов – единица:

$$p_{01} = w_{01} + \min(p_{00} + p_{11}) = 116 + 96 = 212, \quad r_{01} = 1,$$

$$p_{12} = w_{12} + \min(p_{11} + p_{22}) = 58 + 48 = 106, \quad r_{12} = 2,$$

$$p_{23} = w_{23} + \min(p_{22} + p_{33}) = 29 + 24 = 53, \quad r_{23} = 3,$$

$$p_{34} = w_{34} + \min(p_{33} + p_{44}) = 15 + 12 = 27, \quad r_{34} = 4,$$

Разница индексов – двойка:

$$p_{02} = w_{02} + \min(p_{00} + p_{12}, p_{01} + p_{22}) = 142 + \min(64 + 106, 212 + 16) = 142 + 170 = 312, \\ r_{02} = 1,$$

$$p_{13} = w_{13} + \min(p_{11} + p_{23}, p_{12} + p_{33}) = 71 + \min(32 + 53, 106 + 8) = 71 + 85 = 156, \\ r_{13} = 2,$$

$$p_{24} = w_{24} + \min(p_{22} + p_{34}, p_{23} + p_{44}) = 36 + \min(16 + 27, 53 + 4) = 36 + 43 = 79, \\ r_{24} = 3,$$

Разница индексов – тройка:

$$p_{03} = w_{03} + \min(p_{00} + p_{13}, p_{01} + p_{23}, p_{02} + p_{33}) = 155 + \\ + \min(64 + 156, 212 + 53, 312 + 8) = 155 + 220 = 375, \\ r_{03} = 1,$$

$$p_{14} = w_{14} + \min(p_{11} + p_{24}, p_{12} + p_{34}, p_{13} + p_{44}) = 78 + \\ + \min(32 + 79, 106 + 27, 156 + 4) = 78 + 111 = 189, \\ r_{14} = 2,$$

$$p_{24} = w_{24} + \min(p_{22} + p_{34}, p_{23} + p_{44}) = 36 +$$

Разница индексов – четверка:

$$p_{04} = w_{04} + \min(p_{00} + p_{14}, p_{01} + p_{24}, p_{02} + p_{34}, p_{03} + p_{44}) = 162 + \\ + \min(64 + 189, 212 + 79, 312 + 27, 375 + 4) = 162 + 253 = 415, \\ r_{04} = 1,$$

	$p_{11} = 32$	$w_{12} = 32 + 26 = 58$	$w_{13} = 58 + 13 = 71$	$w_{14} = 71 + 7 = 78$
		$p_{22} = 16$	$w_{23} = 16 + 13 = 29$	$w_{24} = 29 + 7 = 36$
			$p_{33} = 8$	$w_{34} = 8 + 7 = 15$
				$p_{44} = 4$

