

Hash Coding の理論と応用

電総研

古川 康一

§ 1. はじめに

hash addressing は、表検索に用いられる技術の1つである。今、情報はキーとデータから成っているものとし、その表を計算機の記憶領域に置いて、キーで検索を行うものとする。情報のとり込み方を、キーと番地とその内容の関係から捉えてみると、図1のような3つの方法が考えられる。a)の index table では、キーのビット系列をそのまま番地として用いる。これは、検索は早いですが、通常大きな記憶領域を必要とする。b)の sequential table では、キーと番地が直接的関係をもたず、或る順序で表をうめる。このとき、キーはデータと同様に記憶させておく必要がある。この方法はa)とは逆に、記憶領域は少なくてすむが、検索が遅くなる。c)の hash table は、a)とb)の中間で、キーの情報的一部分は番地にあり、残りはその内容部にある。この方法は、a)とb)の長所を合わせもっている。

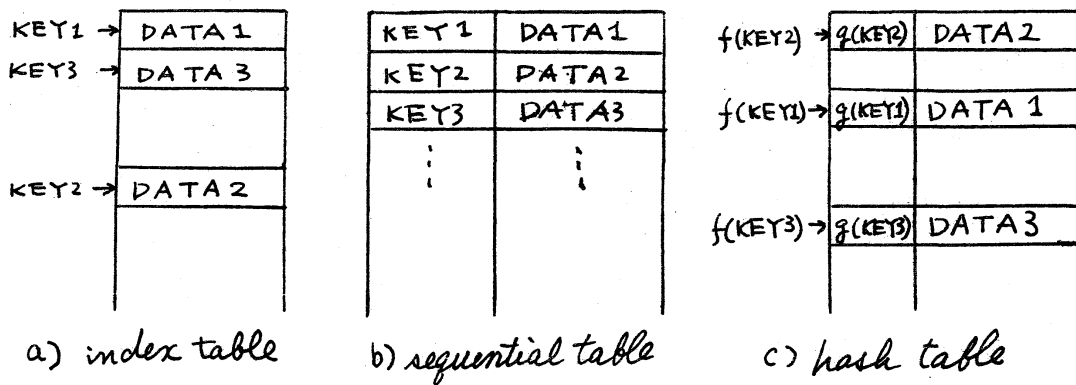


図 1. 各テーブルの構成

§ 2. では, *hash addressing* を実現する変換 (*hash coding*) の具体的な方法, 問題点などを概説する。§ 3. では, *hash addressing* を集合演算に応用したときの改良されたアルゴリズムを与え, その効率の評価を行う。

§ 2. *hash coding* の方法¹⁾

hash coding は, キーの集合から表のエントリのインデックスへの写像と考えることができる。この写像は 1対1 ではないので, 異なるキーが同じインデックスへ写像されることが起こる。この現象をコンフリクト (またはコリジョン) という。良い *hash coding* とは, このコンフリクトが, 統計的に少ない方法である。それは, 各エントリに対する *code* が等確率に発生される場合である。また, キーの集合のもつ規則性を反映しないような方法が良い。これらの点を考慮した代表的な *hash coding* の方法を以下挙げる。

(1) キーのビット系列の一部に定数またはそれ自身をかけ、その結果のビット系列の部分系列を *code* とする。

(2) キーのビット系列の部分系列を加えたものを *code* とする。

(3) キーを表の大きさを割り、余りを *code* とする。

また、キーの集合が与えられている場合には、その集合に依存した写像を適当に見つけ出すことにより、より効率のよい表を作ることができる。²⁾

次に、コンフリクトを解消する方法について述べる。¹⁾ これは、サーチ・アルゴリズムと呼ばれている問題で、大きく2つの方法に分けられる。1つは、コンフリクトを起した項目を同じ表の他の場所におく方法である。他の1つは、コンフリクトを起した項目を、あらかじめ用意された領域において、同じ *hash code* をもつ項目を *linked list* として管理する方法である。この方法は *direct chaining method* と呼ばれ、サーチの効率是他の方法より良い。しかしながら、この方法では余分なメモリ領域を必要とし、自由領域の管理も必要となる。以後の議論は、主に前者について行う。

キー K が与えられたとき、それに相当する表の番地を得る場合に、もし K がすでに表中にあれば、その番地が答となる。もしなければ、表中のまだ使われていない場所を割り当てねばならない。それには、空のエントリが見つかるまで、次々

と、code を発生し続けなければならぬ。今、 i 番目の code を $h_i(K)$ とすると、

$$h_i(K) = h_0(K) + d(i) \quad i \geq 1 \quad (1)$$

と表わせる。ここで $d(i)$ がどのような性質をもてばよいかは、次にあげる 3 つの例によって明らかにされる。

(1) Linear Search $d(i) = i$ または $d(i) = i \times c$ とする方法で、図 3 のようにコンフリクトは必ずクラスタ (cluster) を生じ、そのクラスタ自身は、他のコンフリクトの処理を遅くさせる。

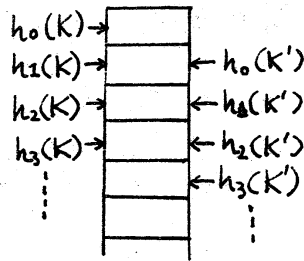


図 3. Linear Search での 1 次クラスタ

(2) Random Search 表の大きさ n 以下の周期をもつ擬似乱数系列を $r_i, i=1, 2, \dots$ とするとき、 $d(i) = r_i$ とする方法で、この方法は、(1) のようなクラスタは生じない。そのため空のセルを見つけるまでの回数 (サーチ回数) は、(1) よりもずっと少なくなる。

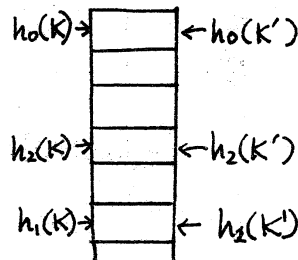


図 4. Random Search での 2 次クラスタ

ただし、一回当りのサーチの手間は、(1) よりもかかる。またこの場合でも、図 4 のような 2 次クラスタを生んでいる。

(3) Quadratic Search³⁾ $d(i) = axi + bxi^2$ とする方法で、(2) に比べて $d(i)$ が簡単に求まる。また、表の大きさが素数ならば、 $d(i)$ は表の半分のエントリに対する code を

発生する。この方法を改良したもので、すべてのエントリに
 対する code を発生するアルゴリズム^{4),5)}、および二次クラスターを
 生じないアルゴリズム⁶⁾が発表されている。

以上の考察から、結論として、良いサーチは、クラスターを生
 じないこと、一回のサーチの手間がかからないこと、記憶領
 域が十分にとれれば Direct chaining method が最も良いこ
 となどがええる。これらの各方法の平均サーチ回数、表の
 占有率 (load factor) に対するグラフを図5に示す。

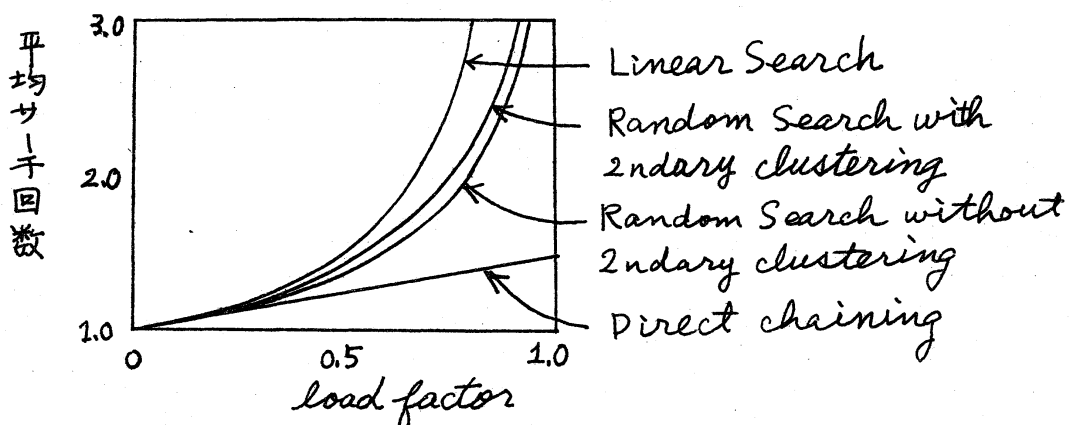


図5. 各サーチの平均サーチ回数

§3. hash addressing の集合演算への応用⁷⁾

キーの集合 A, B を与えて、その合併集合や共通集合を求め
 る等の集合演算に hash addressing を使うことができる。そ
 のとき、基本演算は一方の集合の各要素が他方の集合に含ま
 れているかどうかを知ることである。これは、一方の集合の

すべての要素を hash table に登録して、他方の集合の各要素で検索すればよい。この検索のときに、その表に含まれていない要素に対して、出来るだけ早く含まれていないということがわかればよい。この時間を棄却時間 (reject time) と呼びこことにする。ここでは、わずか1ビットを用いて棄却時間を短くするアルゴリズムを与える。それは、表を作る phase と検索する phase に分け、作る phase でコンフリクトを起したエントリにマークをつけておき (1ビット)、検索するとき、サーチはマークがあるかぎり続けるという方法である。普通、その表に含まれていない要素に対するサーチは空のセルが見つかるまで行われるが、この方法では、マークのない空でないセルにアクセスしてもサーチは終了する。したがって、その分だけ棄却時間が短縮される。

次に、このアルゴリズムによる棄却時間を求める。計算の仮定としては、クラスタは全く生じないものとする。今、hash table の大きさを N とする。 k コの要素を登録したときに、 j コのエントリにマークがつく確率を $P_k(j)$ とする。 $P_k(j)$ に対して、次の漸化式が成り立つ。

$$P_{k+1}(j) = P_k(j) \left\{ \frac{N-k}{N} + \frac{j}{N} \cdot \frac{N-k}{N-1} + \frac{j(j-1)}{N(N-1)} \cdot \frac{N-k}{N-2} + \dots \right. \\ \left. \dots + \frac{j(j-1) \dots \dots \dots 1}{N(N-1) \dots \dots (N-j+1)} \cdot \frac{N-k}{N-j} \right\}$$

$$\begin{aligned}
 & + P_k(j-1) \left\{ \frac{k-j+1}{N} \cdot \frac{N-k}{N-1} + \binom{2}{1} \frac{(j-1)(k-j+1)}{N(N-1)} \cdot \frac{N-k}{N-2} \right. \\
 & \quad + \binom{3}{2} \frac{(j-1)(j-2)(k-j+1)}{N(N-1)(N-2)} \cdot \frac{N-k}{N-3} + \dots \\
 & \quad \left. \dots + \binom{j}{j-1} \frac{(j-1)(j-2)\dots(k-j+1)}{N(N-1)\dots(N-j+2)} \cdot \frac{N-k}{N-j+1} \right\} \\
 & + P_k(j-2) \left\{ \frac{(k-j+2)(k-j+1)}{N(N-1)} \cdot \frac{N-k}{N-2} + \binom{3}{1} \frac{(j-2)(k-j+2)(k-j+1)}{N(N-1)(N-2)} \cdot \frac{N-k}{N-3} \right. \\
 & \quad + \dots + \binom{j}{j-2} \frac{(j-2)(j-3)\dots(k-j+1)}{N(N-1)\dots(N-j+3)} \cdot \frac{N-k}{N-j+2} \left. \right\} \\
 & + \dots \\
 & + P_k(0) \cdot \frac{k(k-1)\dots(k-j+1)}{N(N-1)\dots(N-j+1)} \cdot \frac{N-k}{N-j} \quad (2)
 \end{aligned}$$

この式の $P_k(i)$ の係数は、 k コの要素を登録したとき i コ ($0 \leq i \leq j$) のエントリにマークがあったとき、次の1コの要素を追加してマークが j コになる確率である。その i コ + 1 項は、そのときに i コのすでにマークのあるエントリを途中でアクセスした場合の確率である。

(2)式は、計算によって次式のようになる。(Appendix参照)

$$\begin{aligned}
 P_{k+1}(j) &= \frac{N-k}{(N)_{j+1}} \cdot \left\{ P_k(0) \cdot k(k-1)\dots(k-j+1) \right. \\
 & \quad + P_k(1) \cdot N(k-1)\dots(k-j+1) \\
 & \quad + P_k(2) \cdot N(N-1)(k-2)\dots(k-j+1) \\
 & \quad + \dots \\
 & \quad + P_k(j-1) \cdot N(N-1)\dots(N-j+2)(k-j+1) \\
 & \quad \left. + P_k(j) \cdot N(N-1)\dots(N-j+1) \right\} \quad (3) \\
 j &= 1, 2, \dots, k \quad (N)_{j+1} = N(N-1)\dots(N-j)
 \end{aligned}$$

- オ $j=0$ に対しては

$$P_{k+1}(0) = \frac{N-k}{N} P_k(0) \quad (4)$$

(3) 式より

$$(N-i+1)P_{k+1}(i+1) - (k-i)P_{k+1}(i) = (N-k)P_k(i+1) \quad (5)$$

(5) 式を $i=0, 1, \dots, j$ について加え, (4) 式を代入すると,

$$P_{k+1}(j) = \frac{N-k}{N-j} \left(\sum_{i=0}^j P_k(i) - \sum_{i=0}^{j-1} P_{k+1}(i) \right) \quad (6)$$

ここで, N コのエンタリのうち j コがマークされているときの棄却時間 $t(j)$ は,

$$t(j) = \frac{N-1}{N-j+1} \quad (7)$$

で与えられる。¹⁾ これより, k コの要素が登録されているときの平均棄却時間 \bar{t}_k は,

$$\bar{t}_k = E_k \left(\frac{N-1}{N-j+1} \right) = \sum_{j=0}^{k-1} P_k(j) \cdot \frac{N-1}{N-j+1} \quad (8)$$

で与えられる。 $P_k(j)$ は, (6) 式を用いて数値計算し, \bar{t}_k を求めた。図 6 は, コンフリクト・マークを付けないときの棄却時間との比を示している。ここで, コンフリクト・マークのないときの棄却時間は, 7 式の $t(k)$ で与えられる。

次に, この方法を *direct chaining method* と比較してみる。

今, k コの要素を登録したとき, i コの要素のチェーンのできていたエンタリの数を n_i とすると, 棄却時間 t_d は,

$$t_d = \frac{n_0}{N} + \sum_{i=1}^k \frac{n_i}{N} \cdot i = \frac{n_0}{N} + \frac{k}{N} \quad (9)$$

となる。 n_0 は空のセルの数である。 N と k が十分に大きいと

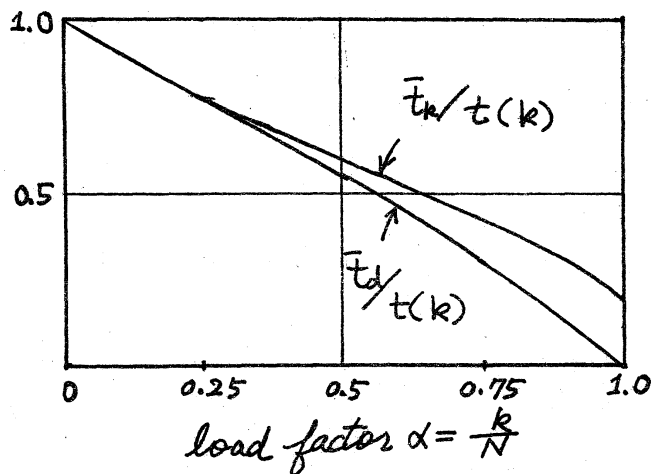
ま、 n_0 は $\lambda = Ne^{-\frac{k}{N}}$ のポアソン分布

$$p(n_0; \lambda) = e^{-\lambda} \frac{\lambda^{n_0}}{n_0!}$$

に従うことが知られている。⁸⁾ n_0 の期待値は、 λ で与えられるので、 $k/N = \alpha$ とすると、(9)の期待値は、

$$\bar{t}_d = e^{-\alpha} + \alpha \quad (10)$$

となる。図6に、 \bar{t}_d と $t(k)$ の比も含めてのせてある。



\bar{t}_k : コンフリクトマークを付けた
ときの棄却時間
 $t(k)$: コンフリクト・マークのない
ときの棄却時間
 \bar{t}_d : direct chaining method
での棄却時間
($N = 1024$)

図6. 棄却時間の比の load factor に対するグラフ

図6. より、load factor が 60~70% になると、コンフリクト・マークを付けたときの棄却時間は、マークのないときの約 ^{半分} になるところがわかる。しかしこの場合にも、direct chaining method がほぼ効率の上限を与えていることがいえる。コンフリクト・マークのゼット数を増やると、この上限に近づくことが予想されるが、どのほうに近づくかは、明らかではない。

参考文献

1. Morris, R. Scatter storage techniques. CACM 11, 1
(Jan. 1968), pp. 38~44
2. 古川, 山本 Hash code に 関 する 命令表の作製法 昭和44年
電気学会連合大会予稿集
3. Maurer, W.D. An improved hash code for scatter storage.
CACM 11, 1 (Jan. 1968), pp. 35~38
4. Radke, C.E. The use of quadratic residue research.
CACM 13, 2 (Feb. 1970), pp. 103~105
5. Day, A.C. Full table quadratic searching for scatter
storage. CACM 13, 8 (Aug. 1970), pp. 481~482
6. Bell, J.R. The quadratic quotient method: A hash code
eliminating secondary clustering. CACM 13, 2 (Feb.
1970), pp. 107~109
7. 古川, 山崎 Hash addressing に 関 する 集合演算. 昭和46
年 電子通信学会 全国大会発表予定
8. Feller, W. An introduction to probability theory and
its applications. Vol I. John Wiley, 1957, pp. 91~95

Appendix. (2)式から(3)式への誘導

$$\begin{aligned}
 f(a, b, c) &= 1 + \frac{b}{a-1} \binom{c+1}{1} \\
 &\quad + \frac{b(b-1)}{(a-1)(a-2)} \binom{c+2}{2} \\
 &\quad + \dots \\
 &\quad + \frac{b(b-1)\dots-1}{(a-1)(a-2)\dots(a-b)} \binom{c+b}{b} \quad (1a)
 \end{aligned}$$

と置いて、この f を用いて (2) 式を書きなおすと、

$$\begin{aligned}
 P_{k+1}(j) &= P_k(j) \cdot f(N, j, 0) \cdot \frac{N-k}{N} \\
 &\quad + P_k(j-1) \cdot f(N-1, j-1, 1) \cdot \frac{N-k}{N} \cdot \frac{k-j+1}{N-1} \\
 &\quad + \dots \\
 &\quad + P_k(0) \cdot f(N-j, 0, j) \cdot \frac{N-k}{N} \cdot \frac{k(k-1)\dots(k-j+1)}{(N-1)\dots(N-j)} \quad (2a)
 \end{aligned}$$

となる。関数 f の一般項 $f(N-j+i, i, j-i)$ を求めると

$$f(N-j+i, i, j-i) = \frac{N(N-1)\dots(N-i+1)}{(N-j+i-1)(N-j+i-2)\dots(N-j)}, \quad (i \geq 1) \quad (3a)$$

となることを、 i についての帰納法を用いて証明することができる。(ただし $i=0$ のときは f は 1 とする)。まず f を書きなおすと、

$$\begin{aligned}
 f(N-j+i, i, j-i) &= 1 + \frac{j-i+1}{N-j+i-1} \binom{i}{1} \\
 &\quad + \frac{(j-i+1)(j-i+2)}{(N-j+i-1)(N-j+i-2)} \binom{i}{2} + \dots \\
 &\quad \dots + \frac{(j-i+1)(j-i+2)\dots j}{(N-j+i-1)(N-j+i-2)\dots(N-j)} \binom{i}{i}
 \end{aligned}$$

となる。 $i=1$ のとき

$$f(N-j+1, 1, j-1) = 1 + \frac{j}{N-j} \binom{1}{1} = \frac{N}{N-j}$$

故に (3a) は成り立つ。

△ $i=h$ のときは (3a) が成り立ったとある。 $i=h+1$ に対して

$$\begin{aligned}
 & f(N-j+h+1, h+1, j-(h+1)) \\
 &= 1 + \frac{j-h}{N-j+h} \binom{h+1}{1} + \frac{(j-h)(j-h+1)}{(N-j+h)(N-j+h-1)} \binom{h+1}{2} + \dots \\
 &+ \frac{(j-h)(j-h+1)\dots(j-h+h-1)}{(N-j+h)(N-j+h-1)\dots(N-j)} \binom{h+1}{h+1} \\
 = & \text{" } \binom{h+1}{l} = \binom{h}{l-1} + \binom{h}{l} \text{ を代入する } \\
 \text{右辺} &= 1 + \frac{j-h}{N-j+h} \left\{ \binom{h}{0} + \binom{h}{1} \right\} + \frac{(j-h)(j-h+1)}{(N-j+h)(N-j+h-1)} \left\{ \binom{h}{1} + \binom{h}{2} \right\} \\
 &+ \dots + \frac{(j-h)(j-h+1)\dots(j-h+h-1)}{(N-j+h)(N-j+h-1)\dots(N-j)} \left\{ \binom{h}{h} + 0 \right\} \\
 &= \frac{j-h}{N-j+h} \left\{ 1 + \frac{j-h+1}{N-j+h-1} \binom{h}{1} + \dots \right\} \\
 &+ 1 + \frac{j-h}{N-j+h} \binom{h}{1} + \frac{(j-h)(j-h+1)}{(N-j+h)(N-j+h-1)} \binom{h}{2} + \dots \\
 &= \frac{j-h}{N-j+h} \cdot f(N-j+h, h, j-h) + f(N-(j-1)+h, h, (j-1)-h) \\
 &= \frac{j-h}{N-j+h} \cdot \frac{N(N-1)\dots(N-h+1)}{(N-j+h-1)\dots(N-j)} + \frac{N(N-1)\dots(N-h+1)}{(N-j+h)\dots(N-j+1)} \\
 &= \frac{N(N-1)\dots(N-h+1)(N-h)}{(N-j+h)(N-j+h-1)\dots(N-j+1)(N-j)}
 \end{aligned}$$

となり, $i=h+1$ に対しても成り立つことがわかった。

この f を (2a) に代入する。(項の順序を逆にする)

$$\begin{aligned}
 P_{k+1}(j) &= P_R(0) \cdot \frac{N-k}{N} \cdot \frac{k(k-1)\dots(k-j+1)}{(N-1)\dots(N-j)} \\
 &+ P_R(1) \cdot \frac{N}{N-j} \cdot \frac{N-k}{N} \cdot \frac{(k-1)\dots(k-j+1)}{(N-1)\dots(N-j+1)} \\
 &+ P_R(2) \cdot \frac{N(N-1)}{(N-j)(N-j+1)} \cdot \frac{N-k}{N} \cdot \frac{(k-2)\dots(k-j+1)}{(N-1)\dots(N-j+2)} \\
 &+ \dots \\
 &+ P_R(j) \cdot \frac{N(N-1)\dots(N-j+1)}{(N-j)(N-j+1)\dots(N-1)} \cdot \frac{N-k}{N}
 \end{aligned}$$

この式を変形すれば, (3) 式が得られる。