

5X-10

Ordered minimal perfect
hashing functionについて

和田 雅美
(東京理科大学 理学部 応用数学科)

1.はじめに

これまで、キー番地変換の一手法として様々なハッシュ関数が提案されてきた[9]。しかし、それらは一般にキー(整数)の集合 $W = \{w_1, w_2, \dots, w_N\}$ から番地の集合 $I = \{0, 1, 2, \dots, M-1\}$ への多対1変換となるため、相異なるキーが同じ番地に変換される現象、いわゆる衝突、を生じるという問題点を持っている。従って衝突が起こった場合の処理についての研究もいろいろと行われて来たが、一方で、衝突の全く起こらない、キーから番地への1対1変換となるハッシュ関数を構成出来ないかという研究が、主に静的および探索と格納を含む準静的ファイルに対して行われてきた。こうしたハッシュ関数を、完全ハッシュ関数(Perfect Hashing Function 以下 PHF)と呼び、その中でもキーの集合の大きさと番地の集合の大きさが等しい、つまりテーブルサイズ最小のものを、最小完全ハッシュ関数(Minimal PHF 以下 MPHF)、その上、データの格納番地の順序がキーの順序に従い、そのシーケンシャル性を保つものを、順配置最小完全ハッシュ関数(Ordered MPHF 以下 OMPHF)と呼ぶ。こうしたPHFの研究の中から、幾つかのPHFおよびその構成法が提案してきた[1, 4, 5, 6, 7]。

Jaeschke [6] は、 $h(w) = \lfloor C/(Dw+E) \rfloor \bmod N$ なる MPHF を提案し、与えられた W に対し、h が MPHF になるような、整数 C、D、そして E を決めるアルゴリズムを与えた(なお、 $\lfloor \cdot \rfloor$ は切り捨てを表す)。井上 [8] は、Jaeschke の関数が OMPHF になるための W に関する必要十分条件を求め、それを用いて、h が OMPHF となるような C、D、そして E の決定アルゴリズムを与えた。Chang [1, 2, 3] は、 $h(w) = C \bmod p(w)$ なる素数関数 p を用いた OMPHF を提案し、与えられた W に対し、h が OMPHF になるように、整数 C を決めるアルゴリズムを与えた。Chang の関数は、Jaeschke の関数に比較し、必ず OMPHF が求まると言う利点があるが、その素数関数の与え方や、整数 C の大きさがかなり大きくなることなど議論の残るところである。以上 2 つのハッシュ関数が、整数論の中中国式剩余定理の応用であるのに対し、Spragnoli [7] は、違った観点から、 $h(w) = \lfloor (w+C)/D \rfloor$ なる PHF を提案し、整数 C、D を決めるアルゴリズムを与えた。Spragnoli のそれは、データの値に偏りがある場合に、最小性が得られないという弱点はあるが、MPHF とならない場合でも、順配置の性質を持っており、(つまり、必ずしも Ordered PHF になる)、また先の 2 つと異なって、新しいデータの追加においてもそれが保持されるという点で、検討する価値はかなり大きいと思われる。

そこで我々は、ここでは、検討の第 1 段階として、Spragnoli の関数が、OMP HF となるための W に関する必要十分条件を考察する。そして、その結果を用いて

OMP HF になる場合の整数 C、D の決定アルゴリズムを与える。また、OMP HF とならない場合に、Spragnoli が導入したカッティングの手法を拡張、適用してみる。

2. 必要十分条件

Spragnoli の関数(Quotient reduction function) $h(w) = \lfloor (w+C)/D \rfloor \quad C \in Z \quad D \in Z^+$ は、キーの集合 $W = \{w_1, w_2, \dots, w_N\} \subset Z$ (以降、考察を容易にするため $w_1 < w_2 < \dots < w_N$ を仮定する) から番地の集合 $I = \{0, 1, 2, \dots, M-1\}$ への関数で、その形から分かるように、データについてのシーケンシャル性を保持している。つまり、 $x \leq y$ に対し、 $h(x) \leq h(y)$ となる。

ここでは先ず、 D_{min} D_{max} を定義し、それを用いて Spragnoli の関数 h が OMP HF になるための、W に関する必要十分条件を与える。

【定義 1】

与えられたキーの集合 $W = \{w_1, w_2, \dots, w_N\}$ に対し、 D_{min} D_{max} を次のように定義する。

$$D_{min} = \max_{j>i} [(w_j - w_i + 1)/(j - i + 1)]$$

$$D_{max} = \min_{j>i+1} \lfloor (w_j - w_i - 1)/(j - i - 1) \rfloor$$

なお、「 $\lfloor \cdot \rfloor$ 」は切り上げを表す。

【定理 2】

与えられたキーの集合 $W = \{w_1, w_2, \dots, w_N\}$ に対し、 $D_{min} \leq D_{max}$ である時、そしてその時のみ、以下のよう C、D を選ぶことで、ハッシュ関数: $h(w) = \lfloor (w+C)/D \rfloor$ は、W について OMP HF となる。

$$D : D_{min} \leq D \leq D_{max} \text{ なる整数},$$

$$C : C_{min} = \max_{j} \{(j-1)*D - w_j\}$$

$$C_{max} = \min_{j} \{(j)*D - w_j - 1\} \text{ によって決まる}$$

$$C_{min} \leq C \leq C_{max} \text{ なる整数}$$

3. アルゴリズム

ここでは、2. で与えられた必要十分条件を用いた OMP HF 構成アルゴリズム(Algorithm OMPHF1)を与える。アルゴリズムは、キー N 個とその個数 N をデータとして与え、それに対して関数 $h(w) = \lfloor (w+C)/D \rfloor$ が OMP HF になるような整定数 C、D の対の 1 つ (C_{min}, D_{min}) を返す(line 18)。データが必要十分条件を満たさない場合は、アルゴリズムはそのメッセージを返して止まる(line 10)。

Algorithm OMPHF1

```

0 procedure OMPHF1(N, w1, w2, ..., wN)
1 if N=1 then return (1, -w1)
2 else begin
```

```

3   Dmin←1; Dmax←∞;
4   for j←2 until N do begin
5     dmin←1; dmax←∞;
6     for i←1 until j-1 do begin
7       dmin←max(Dmin, ⌈(wj-wi+1)/(j-i+1)⌉);
8       if j>i+1 then
9         dmax←min(Dmax, ⌊(wj-wi-1)/(j-i-1)⌋)
10      end;
11     if dmin>Dmax or dmax<Dmin then
12       return ('non-OMPHF')
13     else begin
14       Dmin←max(Dmin, dmin);
15       Dmax←min(Dmax, dmax)
16     end;
17   Cmin←-w1;
18   for j←2 until N do
19     Cmin←max(Cmin, (j-1)*Dmin-wj);
20   return (Dmin,Cmin)
21 end

```

4. カッティング

与えられたキーの集合 $W = \{w_1, w_2, \dots, w_N\}$ においてキーに偏りがあった場合、Sprugnoli の関数 h は、OMP HF からかなり程遠くなってしまう。そこで、彼らは、あるキー w_i を区切り点として、2つのハッシュ関数を接続して、全体でより OMP HF に近づくことを考え、この操作をカッティングと呼んだ[7]。つまり、次のような、整数 $C1, C2$ 、そして D を決定するわけである。

$$h(w) = \lfloor (w+C_1)/D \rfloor, (w_1 \leq w \leq w_j)$$

$$h(w) = \lfloor (w+C_2)/D \rfloor, (w_j < w \leq w_N)$$

そこで、我々はここで、この操作を次のように複数の区切り点に拡張し、全体で OMP HF を構成するようなアルゴリズム (Algorithm OMPHF2) を与える。

$$h(w) = \lfloor (w+C_t)/D_t \rfloor, (a_{t-1} < w \leq a_t), (t = 1, 2, \dots, k)$$

アルゴリズムは、キー N 個とその個数 N をデータとして与え、 (a_t, D_t, C_t) の対を返してくれる。

Algorithm OMPHF2

```

0 procedure OMPHF2(N, w1, w2, ..., wN)
1 if N=1 then return (w1, 1, -w1)
2 else begin
3   t←0;
4   while t<N do begin
5     s←t+2; Dmin←1; Dmax←∞;
6     for j←s until N do begin
7       dmin←1; dmax←∞;
8       for i←s-1 until j-1 do begin
9         dmin←max(dmin, ⌈(wj-wi+1)/(j-i+1)⌉);
10        if j>i+1 then
11          dmax←min(dmax, ⌊(wj-wi-1)/(j-i-1)⌋)
12      end;
13      if dmin>Dmax or dmax<Dmin then begin
14        t←j-1; goto *SKIP
15      end;
16      else begin
17        Dmin←max(Dmin, dmin);
18        Dmax←min(Dmax, dmax)
19      end;
20    end;
21  end;
22  t←N;
23  print (w1, Dmin, Cmin);
24  if t=N-1 then begin
25    t←N; print (wN, 1, N-1-wN)
26  end
27 end
28 end

```

```

20 *SKIP:
21   Cmin←(s-2)*Dmin-ws-1;
22   for j←s until t do
23     Cmin←max(Cmin, (j-1)*Dmin-wj);
24   print (wt, Dmin, Cmin);
25   if t=N-1 then begin
26     t←N; print (wN, 1, N-1-wN)
27   end
28 end

```

5. 例題

1) Sprugnoli の例

$$W = \{17, 138, 173, 294, 306, 472, 540, 551, 618\}$$

$$h(w) = \lfloor (w-2)/76 \rfloor \quad (16 < w \leq 306)$$

$$h(w) = \lfloor (w-287)/37 \rfloor \quad (306 < w \leq 618)$$

2) 12ヶ月の英名の3、2文字目の ASCII コード

$$W = \{20033, 16965, 21057, 21072, 22849, 20053, 19541, 18261, 20549, 21571, 22095, 17221\}$$

$$h(w) = \lfloor (w-16447)/774 \rfloor \quad (16964 < w \leq 20033)$$

$$h(w) = \lfloor (w-17512)/445 \rfloor \quad (20033 < w \leq 22849)$$

6. まとめ

我々は、Sprugnoli が提案した PHF が、OMP HF となるための必要十分条件を求め、カッティングを用いて、OMP HF を構成するアルゴリズムを与えた。

MP HF や OMP HF はあくまで静的で少ないデータに対してのみ構成が実現的ではあるが、一般の動的な大量データ管理に利用できる良いハッシュ関数の選択に対し、それら研究の成果の持つ価値は、大きいものと思われる。

引用文献

- (1) C.C. Chang: The study of an ordered minimal perfect hashing scheme. Comm. ACM 27, 4 (April 1984), 384-387.
- (2) C.C. Chang: The study of a letter oriented minimal perfect hashing scheme. Proc. of the Int. Conf. on Foundations of Data Organization (May 1985), 61-65.
- (3) C.C. Chang and R.C.T. Lee: A letter-oriented minimal perfect hashing scheme. The Comput. J. 29, 3 (June 1986), 277-281.
- (4) R.J. Cichelli: Minimal perfect hash functions made simple. Comm. ACM 23, 1 (Jan. 1980), 17-19.
- (5) S.P. Ghosh: Data Base Organization for Data Management. Academic Press, New York, 1977.
- (6) G. Jaeschke: Reciprocal hashing: A method for generating minimal perfect hashing functions. Comm. ACM 24, 12 (Dec. 1981), 829-833.
- (7) R. Sprugnoli: Perfect hashing functions: A single probe retrieving method for static sets. Comm. ACM 20, 11 (Nov. 1977), 841-850.
- (8) 井上久美子: A method for generating ordered minimal perfect hashing function. 修士論文, 東京理科大学 (March 1986), 1-12.
- (9) 西原清一: ハッシングの技法と応用. 情報処理 21, 9 (Sept. 1980), 980-991.