

高頻度隣接の連結に基づいたデータ圧縮法

The data compression based on concatenation of frequentative code neighbor

中村 博文 *1 村島 定行 *2

Hirofumi NAKAMURA Sadayuki MURASHIMA

Abstract : A parsing algorithm for sequential data compression is given. And a coding algorithm and a decoding algorithm based on this parsing is given. The analysis for parsing is to repeat the operation to regard frequentative neighbor of two code as one new code.

1. まえがき

一次元データの圧縮において、記号列の分解情報を利用する圧縮法がZivとLempelによって提案⁽¹⁾⁽²⁾され発展⁽³⁾⁽⁴⁾⁽⁵⁾⁽⁶⁾されている。いずれも入力データの確率構造を学習しながら圧縮していく適応的な方法である。また、ひとつの切り出しごとに圧縮符号が決まるが、圧縮符号が自己参照的であるため、復号側へ符号表を伝えなくても、復号側では復号しながら符号表を再現できるという特徴を持っている。更に、漸近的最良性が証明されている方法⁽¹⁾⁽²⁾⁽⁵⁾⁽⁶⁾もある。

実用されているLZW法⁽³⁾とその基になった増分分解法⁽²⁾について概要を述べると次の通りである。

増分分解法では、ひとつの切り出しを、過去に登録した（過去に切り出した）もつとも長い一致記号列と1記号とを連結した内容で行い、圧縮符号は一致記号列に登録時につけた番号と1記号の順序数との2つで構成している。LZW法では、ひとつの切り出しを、過去に登録した（切り出し記号列とつぎの記号とを連結

した内容が登録されている）もつとも長い一致記号列で行い、圧縮符号は登録内容につけてある順序数ひとつで構成している。

しかし、以上の方法は次の点が目につく。

(1) 類似の内容の符号表への登録

被圧縮記号列の中に繰り返し現れる記号列がその中間で様々に切り出されて、先頭や末尾が少しずつ異なる記号列が多数登録されている。

(2) 現れない記号列の符号表への登録

登録の時点では使われる可能性が不明なことから、その後現れない記号列も登録している。

これらを次のように解決することを試みた。

(1) 類似の内容の登録の回避

被圧縮記号列中に繰り返し現れる記号列をできるだけひとつの内容として登録する。

そのために先読みをする。ここでは適応的ではなくなるが、ファイルの終わりまで読んで登録と切り出しのための記号列の解析を試みる。2つの記号が隣接して出現する回数が多い隣接から順にひとつの記号と見なし、その内容を基に自己参照的になるように切り出しを行う。

(2) 現れ得る記号列の登録

上記の(1)の先読みした上で出現回数が多い隣接の登録を、出現回数が2以上について行えば、後側で絶対に現れない内容の記号列を登録することはない。

*1 都城工業高等専門学校、都城市
Miyakonojo National College of Technology,
Miyakonojo-shi, 885 Japan

*2 鹿児島大学工学部電子工学科、鹿児島市
Faculty of Engineering, Kagoshima Univ.,
Kagoshima-shi, 980 Japan

2. 解析と符号化

以下4元符号の圧縮例を用いて説明する。

2. 1 使用する用語等

圧縮したい記号列

$$M_0 = | a b a b c b a b a b a a a a a a a a a |$$

があるとする。記号列は、使用記号の集合

$$A_0 = \{a, b, c, d\}$$

の要素の、0回以上の出現による並びである。

A_0 の全要素に $a < b < c < d$ なる大小関係を仮定する。また、更に記号が必要なときには、新たな記号として A_0 の最大の要素に続けて $d < e < f < g$ のような記号をいくらでも使用できるものとする。

2記号が隣合うこと又はこの2記号の組を隣接という。隣接 $a b$ を単に $a b$ と記述する。

$a b$ と $b a$ とは異なる隣接である。

記号列中で2記号の隣接が現れる回数を隣接数という。但し、 M_0 における $a a$ のように同じ記号の隣接を数える場合は、記号列中でひとつの記号を2回は数え上げないことにする。

M_0 で $a a$ の隣接数は3である。

隣接数が最大の隣接を最頻の隣接という。最頻の隣接をなす隣接は複数のこともありうる。

隣接する2つの記号 $a b$ を連結して、ひとつ



の記号とみなすとき、 $a b$ のように記述する。これを2進木とも見なす。'

2. 2 解析

記号列の解析は次のようにして行う。

- (1) 記号列中の最頻の隣接からひとつ選び、記号列中のその隣接をすべて連結する。連結は新たなひとつの記号とみなす。
- (2) (1)を、最大の隣接数が2以上である間繰り返す。

以下例を用いて説明する。

M_0 では隣接数4の隣接 $a b$ と $b a$ がもつとも頻繁に現れている。どちらでもよいがここでは $a b$ を選ぶことにする。

M_0 の $a b$ を連結してひとつの記号にすると次のようになる。

$$M_1 = | \overbrace{a b}^1 \overbrace{a b}^1 \overbrace{c b}^1 \overbrace{a b}^1 \overbrace{a b}^1 a a a a a a a a a |$$

$a b$ を連結してひとつの記号とみなしたことを見出し表に次のように備えておく。

$$T_1 = | (a \ b \ \overbrace{a \ b}^1 \ nil) |$$

nil は $a b$ の連結にまだ新しい記号を割り当てていないことを表す。

次に、隣接数3の隣接 $a a$ を連結して次のようにになる。

$$M_2 = | \overbrace{a b}^1 \overbrace{a b}^1 \overbrace{c b}^1 \overbrace{a b}^1 \overbrace{a b}^1 \overbrace{a a a}^3 \overbrace{a a a}^3 |$$

$$T_2 = | (a \ b \ \overbrace{a \ b}^1 \ nil) (a \ a \ \overbrace{a \ a}^2 \ nil) |$$

更に、隣接数2の $a b a b$ を連結して次のようにになる。

$$M_3 = | \overbrace{a b}^1 \overbrace{a b}^1 \overbrace{c b}^1 \overbrace{a b}^1 \overbrace{a b}^1 \overbrace{a a a}^3 \overbrace{a a a}^3 |$$

$$T_3 = | (a \ b \ \overbrace{a \ b}^1 \ nil) (a \ a \ \overbrace{a \ a}^2 \ nil) |$$

$$(a \ b \ \overbrace{a \ b}^1 \ \overbrace{a \ b \ a \ b}^4 \ nil) |$$

隣接数2以上の隣接はなくなったので、まとめ作業を終える。

2. 3 切り出しと符号化

M_3 、 T_3 の情報を構造を持たない記号の列として符号化する方法は一通りではないがその中でも自己参照的な符号化をひとつを挙げる。

例えば次のように符号化できる。

$$F_1 = 1 1 2 5 3 2 6 1 1 7 7 1 1$$

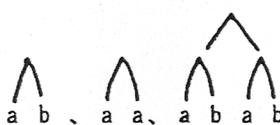
$$F_2 = 1 1 1 0 0 0 0 0 0 1 0 0 0 0 0 1$$

これらをまとめて次のようにも符号化できる。

$$F = 1 1 1 0 1 0 2 0 5 0 3 0 2 0 6 1 0 1 0 1 0 7 0 7 0 1 1$$

F_1 、 F_2 、 F は次のような手順で決定した。

T_3 中のすべてのまとめ方



が、 M_3 に現れる最初の所に印1を付ける。

write(c) : 順序数が c となる記号を出力する
 left(T, c) : 記号 c が表す記号列を構成した左側の記号列の符号を与える
 right(T, c) : 記号 c が表す記号列を構成した右側の記号列の符号を与える
 enter(T, cl, cr, ck, c) : 記号 cl と cr の連結 ck に
 符号 c を割り当てて符号表 T に登録する

```

procedure writestr(c);
begin
  if c<K then write(c);
  else begin
    writestr(left(T,c)); writestr(right(T,c));
  end;
end;
function readtree;
begin
  t:=readtreecode; { from F2 or F }
  if t=0 then begin
    c:=readcode; { from F1 or F }
    writestr(c);
    return(c);
  end else begin
    cl:=readtree; cr:=readtree; Cmax:=Cmax + 1;
    ^~~~~~  

    enter(T,cl,cr,cl,cr,Cmax);
    return(Cmax);
  end;
end;
procedure decode;
begin
  Cmax:=K - 1;
  while not end of file do dummy:=readtree;
end;
  
```

図2. 優化アルゴリズム

4. 評価

ソースプログラム、英文、実行プログラムについて、増分分解法、LZW法と比較した結果を表1. に示す。

LZW法に比べソースファイル等で約2割程度以上、英文と実行ファイルで1割以上圧縮率が向上している。繰り返す内容をひとまとめにして登録し、現れない記号列の登録をなくした効果が現れているといえる。

5. あとがき

以上、被圧縮データを隣接頻度によって解析し、切り出しする方法について述べた。

今後、最後までしている先読みを小量又はゼロにできる適応的なアルゴリズムへの変形をはかりたい。また、処理速度の改善も必要である。

参考文献

(1) J.Ziv and A.Lempel:"A universal algorithm for sequential data compression", IEEE Information Theory, 23,3,pp.337-343(1977).

(2) J.Ziv and A.Lempel:"Compression of individual sequences via variable-rate coding", IEEE Information Theory, 24,5,pp.530-536(1978).

(3) T.A.Welch:"Technique for high-performance data compression", Computer, 17,6, pp.8-19(1984).

(4) 朴志煥, 今井秀樹, 山森丈範, 伊藤秀一, 石井正博: "パターンマッチングと算術符号を用いた実用的なデータ圧縮法", 信学論(A), J71-A, 8, pp.1615-1624(昭63-08).

(5) 曽根和美, 松嶋敏泰, 鈴木謙, 平澤茂一: "ユニバーサルデータ圧縮法の実用化に関する研究", 情報理論とその応用学会, 第11回シンポジウム, pp.79-84(昭63)

(6) 横尾英俊: "Welch型データ圧縮法の新算法", 同上, pp.735-740.

表1. 実験結果

Compression ratio (and entropy assuming i.i.d.)

ファイル名と サイズ[バイト]	増分分解法			LZW法			本方法			圧縮率の改善	
	圧縮率 [%]	切出数 [回]	登録数 [回]	圧縮率 [%]	切出数 [回]	登録数 [回]	圧縮率 [%]	切出数 [回]	登録数 [回]	対増分分解法[%]	対LZW法[%]
lzw.c 10654	83.0 (31.5+22.8+54.4)	2887 (1373)	2887 (1373)	48.2 (44.8)	3727 (1721)	3727 (1721)	32.1+4.3=36.5 (28.1+3.6+31.7)	2895 (1678)	878 (678)	42.1	25.8
parse.c 21084	58.6 (31.4+17.0+46.4)	5096 (2300)	5096 (2300)	47.0 (42.6)	6803 (3034)	6803 (3034)	31.4+4.0=35.5 (27.4+3.3+30.7)	5055 (1209)	1730 (1209)	38.4	24.5
lzw1.c 28028	55.8 (3.06+18.9+47.5)	6350 (3064)	6350 (3064)	44.0 (40.8)	8088 (3849)	8088 (3849)	29.1+3.6=32.8 (25.4+3.1+28.5)	5976 (1460)	2200 (1460)	41.2	25.5
ahuffman.pas 7772	63.1 (32.0+18.4+51.3)	2197 (976)	2197 (976)	50.5 (44.5)	2848 (1267)	2848 (1267)	38.0+4.9=40.8 (30.5+3.9+34.4)	2287 (533)	731 (533)	35.2	18.0
Parse.pas 18481	57.2 (31.3+16.7+46.1)	4417 (2293)	4417 (2293)	44.1 (40.7)	5525 (2828)	5525 (2828)	23.1+3.4=26.5 (18.7+3.1+22.8)	3281 (855)	1675 (855)	53.7	38.8
this.pas 30950	51.0 (28.5+14.2+42.6)	6400 (3259)	6400 (3259)	39.2 (38.6)	7971 (3878)	7971 (3878)	25.0+3.1=28.2 (21.6+2.6+24.2)	5728 (1364)	2025 (1364)	44.7	28.1
chkdsk.exe 10384	78.4 (36.5+30.7+67.2)	3459 (1063)	3459 (1063)	73.6 (60.5)	5216 (1393)	5216 (1393)	57.2+7.0=64.2 (51.5+4.7=56.2)	4773 (775)	1080 (775)	18.1	12.8
diskcopy.exe 19804	76.8 (38.2+28.8+87.0)	6132 (2084)	6132 (2084)	72.0 (60.8)	9102 (2831)	9102 (2831)	54.0+6.4=61.0 (46.3+4.6=52.9)	8021 (1404)	1833 (1404)	20.7	15.3
parse.exe 30304	86.3 (44.0+31.5+75.5)	10252 (3327)	10252 (3327)	82.4 (69.0)	15318 (4670)	15318 (4670)	62.0+6.8=68.7 (54.7+4.7=59.4)	13508 (2308)	3075 (2308)	19.2	15.4
ps.man 11648	57.0 (29.4+15.7+45.0)	2857 (1235)	2857 (1235)	45.2 (40.5)	3744 (1589)	3744 (1589)	33.8+4.3=38.1 (28.3+3.4=31.7)	3075 (714)	943 (714)	33.2	15.7
bs.man 21120	54.7 (29.2+14.7+43.8)	4792 (2101)	4792 (2101)	43.4 (38.2)	6157 (2067)	6157 (2067)	33.3+3.9=37.3 (28.8+2.8=31.6)	5323 (1170)	1332 (1170)	28.5	14.1
sh.man 31232	52.2 (28.8+13.4=42.2)	8803 (3087)	8803 (3087)	40.4 (37.4)	8248 (3878)	8248 (3878)	27.0+3.3=31.3 (24.8+2.5=27.3)	6522 (1485)	1778 (1485)	40.0	21.8

(注1) 圧縮率欄の () 内は理想圧縮率である。増分分解法の左辺は登録符号の順序数と
1符号の順序数を、本方法の左辺は F_1 と F_2 を別々に算出した値である。

(注2) 登録数で 0 内は、1回以上符号生成で使用した登録符号の数で内数である。