

オンライン検索のための文献データ圧縮技法

二村 祥一 松尾 文碩

(九州大学大型計算機センター)

1. まえがき

オンライン文献検索は、大量のディスクスベースを必要とするため、データ圧縮技法を最も必要とする分野のひとつであろう。文献データは抄録、標題など文章テキスト情報の占める割合が大きいので、このデータ圧縮にはテキスト圧縮技法の効率の影響するところが大きい。これまで種々のテキスト圧縮技法が考案され、実用化されている[5]。

テキスト圧縮技法のうち最も本質的で効果的なのは、文字や単語をその生起頻度に応じて、それらを可変長符号で表わす方法である。この方法では、ハフマン符号[3]は最適であることが証明されている[4]。文字要素としたハフマン符号によるテキストの圧縮では、2倍程度の圧縮効果しかなく、それ以上の圧縮効果を上げるには、単語や $N - g r a m$ の文字列を要素とした圧縮方法をとらざるを得ない。しかし、ハフマン符号は、単語のように要素数が非常に大きい場合は、そのまま使えず別の圧縮技法と組み合わせる必要があり、符号化・復号化に要する計算量が大きいという欠点がある。

筆者らは、QOC (Quasi-Optimum Code) と呼ぶ単語を圧縮要素とするテキスト圧縮技法を考案した。QOCは、ハフマン符号のように最適ではないが、ジップの法則が成立する言語に対しては、その圧縮率はハフマン符号より2.7%悪いだけである[1]。一方、符号化・復号化に関する計算量においてQOCはハフマン符号に決定

的に勝る[1]。QOCによりINSPECテープの圧縮を試みたところ抄録、標題のような文章情報のみならず、誌名、著者の所属機関などに対してもこの技法が有効であることが確認できた。すなわち、これらについては約1/4にデータを圧縮することができ、FACOM M-200上で実現した符号器及び復号器の速度は、それぞれ約5.85マイクロ秒/字、1.48マイクロ秒/字と極めて高速であった。

ここでは、QOCを含む文献データの符号化法について述べ、この技法によってINSPECテープを圧縮した結果について報告する。

2. 符号化法

2. 1 QOC

QOCは、基本的には単語を圧縮要素として、その生起頻度の順位(rank)を可変長2進数によって表現した符号である。単語wの順位をrとすると、wの符号 \hat{w} は、4ビットの固定長部分(プレフィックス)とそれに続く可変長の後続部分(マトリックス)からなる。プレフィックスは後続部分の長さを2進数で表示するためのもので、その値は $L_1 \log_2 r$ であり、マトリックスの値は $r - 2^{L_1 \log_2 r}$ である。

自然言語の単語の数は、事実上無限と言ってよいので、単語を要素とする符号化が可能であるのは高順位単語だけであり、低順位単語及び単語の切り出しに用いる複数のデリミ

タに対しては、別の符号化を行わねばならない。QOCでは、プレフィックスの値の0から12までを高順位単語（最高8191語）に、15, 14をそれぞれ低順位単語及びデリミタに用いている。また、小文字をあつかうため13を高順位単語のシフトコードに用いている。

高順位単語のシフトコードは、

- 1) 単語が英大文字のみからなる；
- 2) 単語が英小文字のみからなる；
- 3) 単語の先頭が英大文字で、あとに英小文字が続く；

のいずれであるかを区別するのに用いる。

低順位単語及びデリミタの符号化法には、種々の方式が考えられるが、QOCでは次のような単純な方法を採用した。いま、単語wを $w = a_1 a_2 \dots a_n$ とするとプレフィックス15に続く符号を $\alpha_1 \alpha_2 \dots \alpha_n \#$ とする。 $\alpha_1 \alpha_2 \dots \alpha_n$ は、それぞれ $a_1 a_2 \dots a_n$ の符号で、#はエンドマークである。 $\alpha_1 \alpha_2 \dots \alpha_n$ の中には、必要に応じて文字シフトコード（英大文字、英小文字、数字、その他を区別するためのもの）が挿入される。これらはいずれも5ビット長の符号である。デリミタの符号化は次のように行う。まず、単語間の空白1文字は、単語+空白を符号化したと考え符号化の対象としない。その場合を除いて、おのののデリミタ σ に対して σ , $\sigma\#$, も σ （もは空白）の生起頻度を調べ、その順位に従い可変長2進数0, 100, 101, 11000, 11001, ...をプレフィックス14に続くコードとして割り当てる。可変長2進数の長さは、最初に現れる“0”により判定する。図1にQOCの圧縮表現を、表1, 表2にそれぞれ高順位単語とデリミタ

の符号を示す。

2. 2 1バイト, 2バイトへの符号化

値の異なり数の小さいものについては、おののの値を1バイトあるいは2バイトに符号化する。符号化・複合化には変換プログラム（あるいは変換辞書）が必要になる。

2. 3 4ビット／字, 5ビット／字の圧縮

QOCによる符号化では圧縮効果が少ない書誌事項、あるいは取りうる値の異なり数が大きいものについては、文字単位に4ビット／字あるいは5ビット／字に符号化する。英大文字、英小文字、数字を含む事項について

・高順位単語 ($w = a_1 a_2 \dots a_n$)

$$\hat{w} = \begin{array}{|c|c|} \hline m & r - 2^m \\ \hline \end{array}$$

4ビット mビット

ここでrは単語wの順位。

$$m = \lfloor \log_2 r \rfloor.$$

・低順位単語 ($w = a_1 a_2 \dots a_n$)

$$\hat{w} = \begin{array}{|c|c|c|c|c|} \hline 15 & \alpha_1 & \dots & \alpha_n & \# \\ \hline 4ビット & 5ビット & & & 5ビット \\ \hline \end{array}$$

ここで α_i は a_i のコード。

#はエンドマーク。

・デリミタ (σ)

$$\hat{\sigma} = \begin{array}{|c|c|} \hline 14 & \beta \\ \hline 4ビット & mビット \\ \hline \end{array}$$

ここで β はデリミタ σ の順位による可変長2進コード。

・高順位単語シフト (ϕ)

$$\hat{\phi} = \begin{array}{|c|c|} \hline 13 & \delta \\ \hline 4ビット & 1ビット \\ \hline \end{array}$$

ここで δ の値は0/1で文頭など出現位置により意味が異なる。

図1 QOCの圧縮表現

図るためアセンブラ言語で記述し、IBM 360/370アーキテクチャの性能を十分に引き出すようにしている。

4. INSPECテープの圧縮実験

INSPECテープ [2, 6] 8001～8024（文献数172, 351件）を対象に圧縮実験を行った。INSPECの文献データはA（物理学），B（電気・電子工学），C（計算機・制御工学）の3つの分野に分類できる。INSPECテープ8001～8024での分野別の件数は、A, B, Cそれぞれ108, 141件, 55, 284件, 36, 232件であった。

4. 1 INSPECテープの圧縮方法

抄録、標題、自由索引句の圧縮にはQOCを適用する。QOCでは単語の切り出しのためのデリミタを指定できるが、ここでは英数字を除くすべての特殊記号を登録する。

誌名、著者の所属機関などの圧縮にはQOCを適用する。ただし、誌名、著者の所属機関などの場合は、抄録などのような文章とは違って、使用される単語が限られ、又省略形が多く用いられている。このため抄録、標題、自由索引句とは別に扱う。デリミタとして“ち”，“，”，“-”，“'”，“/”，“\$”の6文字を登録する。次の書誌事項もこのクラスに含める。

会議名、会議開催場所、出版社名、

出版社所在地、後援機関

言語名は、言語単位に1バイトに符号化する。一次文献の記載言語名の異なり数は50以下である。言語名では一つの論文が複数の言語を用いて記載されている場合があり、それを区別する必要がある。

その他の書誌事項は文字単位に4ビット／字あるいは5ビット／字に符号化する。書籍番号のように使用される文字が、数字やいくつかの特殊記号に限られている場合は、4ビット／字に符号化し、それ以外は5ビット／字に符号化する。

- ・ 4ビット／字に符号化する書誌事項
　　書籍番号、分冊番号、ページ数
- ・ 5ビット／字に符号化する書誌事項
　　著者名、CODEN、レポート番号、雑誌のボリューム番号、参照ページ、発行年月日

ここで、雑誌のボリューム番号では、“Vol. ” “Ser. ”, “No. ”が、発行年月日では“Jan. ”, “Feb. ”, “”が固定的に使用されているので、これらについては、文字列全体を5ビット／字に符号化する。著者名については、ほとんどの場合“Smith, A. B.”のようにファミリー名“Smith”に続いてイニシャル“A. B.”が設定されており、これを配慮して符号化した。

4. 2 辞書の作成

本実験では、符号化・復号化のために次の2種類の辞書を作成した。

- 1) 抄録、標題、自由索引句のための高順位単語辞書とデリミタ辞書
- 2) 誌名、著者の所属機関などのための高順位単語辞書とデリミタ辞書

辞書1は分野別に、辞書2は全分野のデータを用いて作成した。辞書1を分野別に作成したのは、分野により使用される単語が異なるためであり、辞書2で全分野のデータを用いたのは、誌名、著者の所属機関などにおいて

合、生起頻度では、それぞれ24%, 31%, 40%; 50%を占めるにすぎない。高位の4096を辞書に乗せ圧縮率の改善を試みたところ10%程度改善されるにすぎない。このため、現在ファミリー名の辞書化は行っていない。著者名ではファミリー名の部分文字列への分解も考えられる。

5. あとがき

ここでは、QOCと呼ぶ単語を圧縮要素とする符号化法とこれを中心とした2次文献データの圧縮法について述べ、続いてINSPECテープを対象に圧縮実験を行い、この技法の性能を評価した。この実験では、文章情報の圧縮率は4以上、2次文献データ全体では圧縮率は3.67という非常に良好な結果を得た。また、オンライン文献検索などへの応用において重要な復号化時間がFACOM M-200で1文字当たり1.68μs程度であり、INSPECのデータをCPU時間1秒で760件程度を復号することができたので、この技法が実用に耐えることがわかった。

九州大学大型計算機センターでは、現在、富士通会話型情報検索システムFAIRS-Iを用いてINSPECの検索サービスを行っている。我々は、AIRと呼ぶ効率を重視した情報検索システムを開発しFAIRS-Iと置き代える計画である。AIRでは、ここで述べた圧縮技法を用いる。

参考文献

- [1] 松尾文碩、二村祥一、吉田将：順最適テキスト圧縮符号、九大工学集報、55, 2, 103-106 (1982).
- [2] 松尾文碩、二村祥一、吉田将：科学技術論文抄録における単語の統計的性質、同上、54, 4, 411-416 (1981).
- [3] Huffman, A.D.:A Method for the Construction of Minimum Redundancy Codes, Proc.IRE, 40, 9, 1098-1101 (1952).
- [4] Gilbert, E.N. and Moore, E.F.: Variable-length Binary Encodings, Bell Syst. Tech. J., 38, 4, 933-967 (1959).
- [5] Radue, J.E.:Text Compression Techniques, Quaestiones Informaticae, 1, 1, 30-36(1979).
- [6] INSPEC Tape Services Manual, The Institution of Electrical Engineers.