

4E-4

## 重ね合せ符号を利用する 大規模辞書検索システムについて

玉越 靖司  
松下電器産業(株) 東京研究所

### 1. はじめに

我々は自然言語処理システム用の大規模な単語辞書と、それに対して辞書引きを行うシステムの開発を進めている。このシステムの目的は、各々の自然言語処理システムで個別に行われている辞書や辞書引きシステムの開発に係わるオーバーヘッドを吸収することである。従って、このシステムにおいて対象とする辞書はマスター辞書の性格をもった汎用大規模辞書である。現在、20万語規模の日本語単語辞書を念頭に置いている。本稿では、重ね合せ符号を利用する辞書引き手法を提案する。

### 2. 理想とする辞書引き

現在、自然言語処理システムで行われている辞書引きとして最も一般的な方法は転置ファイルを用いる方法である。この方法は大規模な辞書に対しても高速に検索が可能である。転置ファイル方式にハッシングを導入して、我々が目指す規模の辞書を非常に高速に検索する方式もある。しかし、こうした転置ファイルを用いる方法では、特に辞書引きのようにマルチキーの検索を行いたい場合、必要な記憶領域量や辞書のメンテナンスのコストが大きい。

一方、例えばPrologのような言語で記述している自然言語処理システムにおいては、述語で辞書を記述し、検索を言語の单一化機能に委ねている例もある。こうした方法は辞書の規模が小さな場合には有効だが、我々が目指す大規模辞書の検索方式としては適当でない。

以上を考え合わせると、我々が理想とすべき辞書引きシステムの特長は次のようなものである。

- (1) 高速な検索が可能である
- (2) 大規模辞書の検索が可能である
- (3) マルチキーの検索が可能である
- (4) 記憶領域量を抑えることが可能である
- (5) メンテナンスが容易である

さらに、種々の自然言語処理システムで用いることを考慮すると、それぞれの環境的制約に対応できるように、これらの条件のトレードオフによる設計が容易に可能であることが望ましい。

### 3. 重ね合せ符号を利用するデータ検索手法

重ね合せ符号を利用するデータ検索の手法(以下SCW方式と呼ぶ)[1]について述べる。SCW方式は、本質的にはパンチカードによるハンドソーティングの手法と同じである。これまで、電話番号案内[2]や文献情報検索[3, 4, 5]、ホーン節の单一化可能項検索[6, 7,

8, 9]などに応用してきた。

#### 3.1 データの蓄積

マスターファイルF中のレコードR<sub>i</sub>が持つ#R<sub>i</sub>個のキーワード1つ1つに対し、ハッシュ関数を用いてb c w(Binary Code Word)と呼ぶ長さbのビット列を作る。このハッシュ関数はb c wのbビットのうち、wビットに'1'を立て、残りは'0'とする。wを重みと呼ぶ。#R<sub>i</sub>本のb c wをビット毎に論理和をとった長さbのビット列s(R<sub>i</sub>)をR<sub>i</sub>のs c w(Superimposed Code Word)と呼ぶ。こうして作成された#F本のs c wを、b本の長さ#Fのビット列としてファイルS C Wに格納する。ここで、#FはF中のレコードの個数である。

このようにして作られるS C Wは、Fの転置ファイルに比べてはるかに小さく、メンテナンスも容易である。しかも、bを決めることにより、S C Wのサイズを自由に設計することができる。

#### 3.2 検索

質問Qが持つ#Q個のキーワードについて、蓄積の場合と同じハッシュ関数を用いてb c wを作成し、同様の処理によりマスクs(Q)を作る。s(Q)とS C Wを照合して、ドロップと呼ばれる検索結果

$D(s(Q)) = \{R_i \mid s(R_i) \wedge s(Q) = s(Q)\}$   
を求める。s(Q)で'1'が立っている位置に対応するS C Wの長さ#Fの列に対して、ビット毎の論理積をとることによってDが簡単に求められる。この論理積の結果である長さ#Fのビット列において'1'が立っているビットiに対応するR<sub>i</sub>がDに含まれる。SCW方式では、#Fのサイズにはほとんど依存せず、高速にDの検索が可能である。

#### 3.3 フォルスドロップ

Qに対する正しい検索結果のレコード集合をA(Q)とすると、D(s(Q))は次の性質を持っている。

$$D(s(Q)) = A(Q) \cup Df(s(Q)) \\ (A(Q) \cap Df(s(Q))) = \emptyset$$

ここで、Dfは誤って検索されたレコードの集合のことであり、フォルスドロップと呼ばれる。SCW方式ではDからDfを除去するオーバーヘッドが必要である。この処理時間コストのオーダーは#Dであり、Dの検索に比べてずっと大きい。ここで#DはD中のレコードの個数である。しかし、bやw、ハッシュ関数の設計により、Dfを十分に小さくすることが可能である。

#### 4. SCW方式の辞書引きへの応用

SCW方式を自然言語処理の辞書引きに応用することを提案する。

##### 4.1 対象とする辞書

まず、辞書引きの対象となる辞書の性質を規定する。我々は20万語規模の日本語単語辞書で「見出し語情報」「品詞情報」「意味情報」を持っているマスター辞書[10]を想定している。

これらの情報は階層構造になっている。すなわち、1つの「見出し語情報」は1つ以上の「品詞情報」を持っている。また、1組の「見出し語情報」「品詞情報」は1つ以上の「意味情報」を持っている。

この辞書引きシステムを用いる自然言語処理システムとして、日本語解析と日本語生成の両方を考慮する。日本語解析システムでは、主に「見出し語情報」と「品詞情報」をキーとする辞書引きが行われる。日本語生成システムでは、主に「意味情報」をキーとする辞書引きが行われる。

##### 4.2 SCW方式の辞書引きの特徴

SCW方式が持つ一般的な性質を以下にまとめる。

- (1) #Fにはほとんど依存せず、高速にDが求まる
- (2) 検索時間コストの大部分を占めるDf除去は、オーダー#Dで抑えられる
- (3) bを調節することにより、SCWのサイズを決定できる
- (4) bを大きくすると、SCWは大きくなるが、#Dfは小さくなり検索が高速化される
- (5) s(Ri)中の'1'の平均個数がb/2になるようwを調節することにより、#Dfが最小となる
- (6) マルチキーの検索が容易で、しかも#Qが大きいほど#Dが小さくなり、検索が効率的である
- (7) Riに関するメンテナンスはs(Ri)だけによく、簡単である

こうした特徴は、本稿で理想とする辞書引きに適合しており、SCW方式が辞書引き方式として有望であることを示している。

しかし、辞書が持つキーワードのb c wを単純に重ねさせるだけでは、十分に#Dfを小さくすることができない。そこでSCW方式の拡張として、キーワードの性質により、scw中でのb c wの位置(フィールド)を分ける方法を探る。

##### 4.3 見出し語情報のフィールド

SCW方式では、それぞれの異表記に対応するb c wを作成することにより、容易に異表記の辞書引きを実現することが可能である。また、見出し語の文字列の前方文字から順にキーワード化することにより、前方一致の辞書引きも可能である。

これらの機能を実現するためには、見出し語情報のフィールドにおいて多重のb c wの重ね合せを行う必要がある。従って、見出し語情報のフィールド幅bmは他のフィールド幅に比べてかなり大きくしなければならない。

#### 4.4 品詞情報のフィールド

品詞情報では、品詞や活用など、キーの種類の数が限られている。また、b c wの重ね合せの数が決まっている。一般に、重ね合せが少ないとscw中の'1'が少なくなり、Dfの発生が抑えられる。逆に言うと、bに対するwを相対的に大きくできることによって、Dの精度が上がる。

こうした特徴を利用して、品詞情報のフィールド幅bhはbmに比べて小さく設定することが可能である。

#### 4.5 意味情報のフィールド

意味情報によるキーは一意の概念番号である。従って、意味情報のフィールドではb c wの重ね合せが行われず、他のレコードとの重複によるDfがない。このことは、#F = 20万のとき意味情報のフィールド幅bi = 21とすれば、Dfを全く発生させないことを意味する。

#### 4.6 頻度別SCWによる高速化

辞書引き全体の効率のためには、使用頻度の高い見出し語を特に高速に引く必要がある。そこで、#Dを小さくすることによってオーダー#Dの処理時間コストがかかるDf除去のオーバーヘッドを小さくする。すなわち、辞書中の単語頻度情報を利用して頻度別SCWを作り、高頻度のSCWから順に検索を行う。

#### 4.7 辞書の定量的解析による準最適化

スタティックな辞書に対して定量的な解析を行うことにより、またこの辞書引きシステムを用いる自然言語処理システムの特性に応じて、各々のフィールドで準最適なbやw、ハッシュ関数、頻度別SCWのサイズを設計することが可能である。

### 5. おわりに

SCW方式は処理時間や記憶領域量、メンテナンスのコストに優れ、環境的制約によるこれらのトレードオフが可能であるということに着目して、SCW方式による辞書引きシステムを提案した。SCW方式ではDfの除去というオーバーヘッドがあるが、その軽減が十分に可能であるとの見通しを得た。現在、この方法の有効性を確認するために、実験システムを作成している。準最適化の手法とその結果については次の機会に報告したい。

#### 参考文献

- [1] Roberts, C. S. : Partial-Match Retrieval via the Method of Superimposed Codes, Proceedings of the IEEE, Vol. 67, No. 12, 1979.
- [2] Gabbe, J. D. et al. : Applications of Superimposed Coding to Partial-Match Retrieval, Proc. COMPSAC 78, 1978.
- [3] 有川他：重ね合せ符号と逐字サーチを利用する文献情報検索システムについて(1), 情報処理学会第25回全国大会, 1982.
- [4] 有川他：重ね合せ符号と逐字サーチを利用する文献情報検索システムについて(2), 情報処理学会第27回全国大会, 1983.
- [5] 有川他：重ね合せ符号と逐字サーチを利用する文献情報検索システム, 情報学シンポジウム講演論文集, 1985.
- [6] Wise, M. J. et al. : Indexing PROLOG Clauses via Superimposed Code Words and Field Encoded Words, Proc. 2nd Int. Logic Progr. Conf., 1984.
- [7] 森田他：スーパーインボーズドコードを用いた構造体の検索方式, 情報処理学会第33回全国大会, 1986.
- [8] 和田他：KBMS PHI:重ね合せ符号を用いた関係演算の処理方式, 情報処理学会第34回全国大会, 1987.
- [9] 森田他：MPPMを用いた知識ベースマシン(3)-構造体のインデックス方式に関する考察-, 情報処理学会第35回全国大会, 1987.
- [10] 単語辞書(第2版), EDR TR-006, 1988.