

スーパーインポーズドコードを用いた構造体の検索方式

森田 幸伯 和田 光教* 伊藤 英則

(財) 新世代コンピュータ技術開発機構 * (株) 沖電気工業

1. はじめに

最近エキスパートシステムや知識ベースなどでは知識を表現するものとして、変数を持った構造体が多く使用されている。このように表現された知識の扱いには、パターンマッチ、特に单一化処理が重要な役割を演じている。今後、扱う知識の量が増大すると併に、1次記憶や2次記憶における知識の高速な検索技術がますます重要となってくる。

一方重ね合わせ符号(superimposed code)法を用いた検索技術は、マルチキーファイルの構成やシグニチャファイルを用いたファイルアクセスの高速技法などに適用されている(1)。また、エディタなどのサードコマンド等にも応用されてきている。

本稿では、重ね合わせ符号法を構造体の一項である項に対して適用させる方式について述べる。これは、ある項と单一化可能な項を含むオブジェクト(レコード、ページ等)に対するインデックス検索を可能にするものである。

2. 定数に対するSCW方式(1)

重ね合わせ符号による検索方式(本稿ではSCW方式とよぶ)は、レコードに対してsuperimposed code word(SCWと記す)なる一種のインデックスを用いて、大容量のレコードを検索する時間を短縮しようとするものである。

R_i をファイルFのレコードとする。 R_i には幾つかのキーワードが付けられている。各キーワードはある関数(ハッシュ関数)によって2進符号語(BCW)に写像される。あるレコード R_i に対してその全てのキーワードのBCWのビットごとのORをとったものをSCWと呼び、 S_i と表記する。ファイルFの全てのレコードについて対応するSCWを計算し1つのファイルとし、これをSCWファイルSと呼ぶ。

今、ファイルFに対して幾つかの指定したキーワードをもつレコードを検索する問題を考える。この検索のために指定されたキーワードを上記SCWを生成したのと同じ関数を用いてBCWを生成しORをとってSCWを生成する。これをキュアリマスクと呼びQで表わす。QとあるSCW・ S_i とが以下の条件(1)を満足しないならば、対応するレコード R_i は、はじめの検索条件(キュアリ)を満足しない。

$$Q \wedge S_i = Q \quad (1)$$

ファイルSから(1)を満足する S_i の集合を選び出し、対応する R_i に対してのみキュアリを満足するかどうか調べることにより探索空間を減らすことができる。勿論条件(1)を調べる処理は、オーバーヘッドとなる。しかし、この処理は比較的簡単な処理であるため、SCWをうまく設計すれば、このオーバーヘッドは、十分小さくできる。

3. SCW方式を用いた項の検索方式

前節のキーワードを変数を含む一種の構造体である項に拡張する。それにより、キュアリはある項と单一化可能な項をキーワードとしてもつレコード(オブジェクト)を求めるものとなる。(本稿では議論を簡単にするために各レコードのキーワードもキュアリのキーワードも1つとする。)

- A. 項が構造体である(木構造とみなせる)ことより、木構造の各ノードに対してそのノードの値(関数記号または変数記号)から2進符号語の一部分への関数(特定ビットを0に固定する)を用意し、項のBCWとしては、それらの値のORをとったものとする。
- B. 項が変数を含むことにより、キーワードからBCWを生成する関数を通常のデータに対する関数とキュアリに対する関数を別の関数とする。

結局、項 t にたいするBCWは、 t を木構造とみなしたときのノードの集合をNodeとし、ノード n のデータ用関数を h_n 、ノードにある値(関数記号or変数記号)をname_nとすると

$$BCW_t = h_t(t) = \bigvee_{n \in Node} h_n(name_n) \quad (2)$$

となる。同様にキュアリに現われる項 t にたいしては、ノード n のキュアリ用関数を h_n' とすると、

$$BCW_t' = h_t'(t) = \bigvee_{n \in Node} h_n'(name_n) \quad (3)$$

となる。

各ノードのデータ用関数またはキュアリ用関数 h_n に対し、全ての関数記号および変数記号 x に対してある $\{a_i\}_{i=0,1}$ が存在して、 $h_n(x) = \sum_{i=0}^k a_i 2^{i-1}$ となる最小の整数の集合EをBIT(h_n)で表わす。

各のノードのデータ用関数およびキュアリ用関数に対し

Structure retrieval via the method of superimposed codes
Yukihiro MORITA¹, Mitsuomi WADA², Hidenori ITOH¹
¹ICOT, ²Oki Electric Industry Co., Ltd.

て以下の条件を考える。

【条件1】

親ノード n を定めると、親ノードのデータ用およびキュアリ用の2進符号への関数 h_n, h_n' に対し、 m 個の子ノードのデータ用およびキュアリ用の2進符号への関数 $h_{n_i}, h_{n'_i}$ ($1 \leq i \leq m$) がそれぞれ一意に定まり、

$$\begin{aligned} BIT(h_n) &\subseteq BIT(h_{n'}) \\ BIT(h_{n'_i}) &\subseteq BIT(h_{n'}) \end{aligned}$$

【条件2】

全てのノードの2進符号への関数は、以下を満足する。任意の変数記号 x に対して、 $E = BIT(h_n)$ とすると、

$$\begin{aligned} h_n(x) &= \sum_{i \in E} 2^{i-1} \\ h_{n'}(x) &= 0 \end{aligned}$$

任意の関数記号 x に対して、

$$h_n(x) = h_{n'}(x)$$

【補題】

条件1、2を満足する項 t のデータ用ハッシュ関数 $ht(t)$ 、及びキュアリ用ハッシュ関数 $ht'(t)$ は、

$$\begin{aligned} ht'(t_1) \wedge ht(t_2) &= ht'(t_1) \\ ht'(t_2) \wedge ht(t_1) &= ht'(t_2) \end{aligned}$$

をみたす。ただし、 θ は置換(substitution)。

【定理】

条件1、2を満足する項 t のデータ用ハッシュ関数 $ht(t)$ 、及びキュアリ用ハッシュ関数 $ht'(t)$ は、 t_1, t_2 が単一化可能ならば

$$ht'(t_1) \wedge ht(t_2) = ht'(t_1) \quad (2)$$

をみたす。

【証明】

t_1, t_2 が単一化可能なことよりある置換 θ が存在して $t_1\theta = t_2\theta$ かつこれらの項が変数を含まないようにすることができる。

補題より

$$\begin{aligned} ht'(t_1) \wedge ht(t_2) &= ht'(t_1) \\ ht'(t_1) \wedge ht(t_2) &= ht'(t_1\theta) \\ t_1\theta \text{ に変数を含まない} \text{ ことから,} \\ ht(t_1\theta) &= ht'(t_1\theta) \end{aligned}$$

これらより

$$\begin{aligned} ht'(t_1) \wedge ht(t_2) &= ht'(t_1) \\ ht'(t_1) \wedge ht'(t_2) &= ht'(t_1) \\ ht'(t_1) \wedge ht'(t_2\theta) \wedge ht(t_2) &= ht'(t_1) \\ ht'(t_1) \wedge ht(t_2\theta) \wedge ht(t_2) &= ht'(t_1) \\ ht'(t_1) \wedge ht(t_2) &= ht'(t_1) \quad \square \end{aligned}$$

本稿では定理2の(2)を満足する t_2 を t_2 の疑似単一化可能項と呼ぶ。

4. 簡単な評価と他方式との比較

单一化処理をおこなう前にその可能性のある項の組を選び出す方式は、変数のバインディングを無視する方法[2]、項を文字列で表わし変数の前までを比較する方法[3]、構造を限定しハッシュ関数を用いて符号化して調べる方法[4]等幾つか提案されている。

表1は、幾つかのサンプルデータに対する疑似单一化項と单一化可能項の割合を示している。BCWは26ビットとし、条件1については n 引数関数に対して BCWのビットを $n+1$ 等分して各子ノードに割り当てる。データは1引数のものから4引数のものまで用意し、各引数は3種の定数または変数である。

これらの例はハッシュ関数に衝突の無い場合であり[2]や[4]の方式と同程度の割合であると思われる。[3]に対しては平均7割位低くなっている。衝突が多くなってくると[2], [4]よりも单一化可能でない疑似单一化項を選び出す確率が高くなると思われる。しかし、本方式は[2]よりも簡単な処理(ビット演算のみ)で疑似单一化可能項を選択でき、[4]のように構造に対して制限を必要としない。また、複数の項のどれかと单一化可能であるという検索や複数のキュアリに対しては有効であると思われる。

表1. 疑似单一化可能項と单一化可能項の比率

| 引数の数 | 1 | 2 | 3 | 4 |
|-------------|-------|-------|-------|-------|
| 疑似单一化可能項の比率 | 1.000 | 1.094 | 1.317 | 1.667 |

5. おわりに

重ね合わせ符号法が項に対する单一化可能条件に応用できることを示した。3節ではレコード(オブジェクト)とキュアリに含まれる項の数を1つとしたが、通常のSCWと同様に3節の項に対するBCWを重ね合わせて使用することも出来る。また、本方式の有効な具体的な応用例およびそのときのハッシュ関数の遊びかた、特に条件1の各ノードの関数の定めかた等も今後の課題である。

【参考文献】

- [1] Roberts, C. S., "Partial-Match Retrieval via the Method of superimposed Codes", Proc. of IEEE, Vol.67, No.12, pp.1624-1642, Dec 1979.
- [2] Sabbarel, G. B., "Unification for A Prolog Database Machine" 2nd Logic Programming Conference 1984.
- [3] Morita, Y. et al."Retrieval-By-Unification Operation on a Relational Knowledge Base Model", Proc. of the 12th Int'l Conference on VLDB 1986.
- [4] 大森 他「推論機能と関係データベースの融合－ハッシュとソートによる検索」第32回情報処理大, 1M-6, 1986.