

シグネチャを持つ入れ子型索引を用いた集合値検索のコスト解析

永戸 克裕[†]北川 博之^{††}[†]筑波大学 理工学研究科 ^{††}筑波大学 電子・情報工学系

1 はじめに

集合は基本的なデータ構造であり、複合オブジェクトの部分構造としても頻りに現れる。このため、集合値に関する検索条件を効率的に支援する索引機構は、先進的な応用分野を対象としたデータベースシステムにおいて重要なものとなる。

集合値検索を支援する索引機構として、入れ子型索引(NIX)や、シグネチャファイルなどがある。この2つの索引機構を比較すると、問合せ条件やデータベースの規模によってその優劣は異なる。そこで、シグネチャとB+木を組み合わせることで、それぞれの長所を生かしたより優れた集合値検索のための索引機構を提案し、そのコストを評価する。

2 集合値検索とその処理

2.1 集合値検索

例として学生に関するデータベースを考える。Studentクラスには属性としてnameとhobbiesがあり、hobbies属性は文字列の集合を値とする。ここで、次の問合せ Q_1 を考える。

```
Q1: select name
    from Student
    where hobbies ⊇ {"Baseball", "Fishing"}
```

Q_1 は、趣味に“Baseball”と“Fishing”の両方を持つ学生の名前を求める問合せである。問合せ条件中に現れる集合を問合せ集合(query set, Q)と呼び、データベース中に格納され問合せ集合との比較の対象となるそれぞれの集合をターゲット集合(target set, T)と呼ぶ。

問合せ集合を部分集合として含むターゲット集合を検索する問合せを $T \supseteq Q$ (has-subset)の問合せと呼び、問合せ集合の部分集合となるターゲット集合を検索する問合せを $T \subseteq Q$ (is-subset)の問合せと呼ぶ。

2.2 集合値検索の処理

集合値検索を支援する索引機構として、入れ子型索引(NIX)とシグネチャファイルについて考える。

NIXは、B+木に基づく索引手法で、リーフページにキー値と対応する識別子(OID)リストからなる索引エントリが格納される。 $T \supseteq Q$ の問合せでは、与えられた D_q 個の要素を持つ問合せ集合に対して、対応する D_q 個のOID集合が検索される。次に、このOID集合の積集合中の各OIDを基にデータオブジェクトが検索される。 $T \subseteq Q$ の場合は、同様に D_q 個のOID集合が検索され、次にそれらの和集合を求める。この和集合に対応するデータオブジェクトは必ずしも検索条件を満たすとは限らないため、実際にデータオブジェクトを検索して、条件を満たすか調べなければならない。

シグネチャファイル(signature file)は、主にテキストデータベースにおけるキーワード検索において利用されてきた索引手法である。シグネチャとは、個々のデータオブジェクトから生成される固定長のビット列であり、

シグネチャを対応するデータオブジェクトのOIDとともに格納したものがシグネチャファイルである。

集合シグネチャは以下のように生成される。集合値が与えられると、まず、ハッシュ法などにより、集合の各要素から要素シグネチャ(element signature)が作られる。次に、すべての要素シグネチャのビットごとの論理和をとることにより、集合シグネチャ(set signature)が生成される。作成された集合シグネチャ(ターゲットシグネチャ(target signature))は、対応するOIDとともにシグネチャファイルに格納される(図1)。

集合の要素	要素シグネチャ ($F=8, m=2$)
“Baseball”	→ 01000100
“Golf”	→ 00100001
“Fishing”	→ 00010100
	↓ 論理和
集合シグネチャ	01110101

図1: 集合シグネチャの生成

シグネチャファイルの物理的な構成手法として、ビットスライスシグネチャファイル(Bit-sliced Signature File, BSSF)[1, 2]がある。BSSFではシグネチャはビットごとに別々のファイル(ビットスライス)に格納されるため、検索では一部にアクセスすばよいことになる。

$T \supseteq Q$ に対する問合せ処理は以下ようになる。問合せが与えられると、問合せ集合より問合せシグネチャ(query signature)が生成される。次に、問合せシグネチャにおいて“1”の値を持つ各ビット位置に対応するビットスライスが検索される。次に、読み出したすべてのビットスライスの論理積をとり、結果“1”の値を持つエントリについて、対応するOIDをOIDファイルから検索する。対応するデータオブジェクトは問合せを満たす候補(ドロップ(drop))となる。最後に、それぞれのドロップが実際に問合せ条件を満たすか否か調べる。

$T \subseteq Q$ に対する問合せ処理は、 $T \supseteq Q$ の処理で“1”と“0”の役割を入れ換えたものである。

検索対象のデータベースが小・中規模の場合、 $T \supseteq Q$ ではNIXとBSSFは同程度の性能であり、 $T \subseteq Q$ ではBSSFの方が優れる。しかし、データベース中のオブジェクト数が増加すると、BSSFに対するNIXの優位性は徐々に低減する。

3 シグネチャとB+木を組み合わせた索引

上の2つの索引機構の長所を生かすため、次のようにしてこの2つを統合することを考える。

シグネチャ付きNIX(S-NIX)は、NIXのリーフページに、キー値だけでなく、対応するOIDリスト、ターゲットシグネチャが格納される(図2)。

3.1 検索処理

S-NIXにおける、 $T \supseteq Q$ に対する検索処理はNIXと同様で、シグネチャ部は使用しない。 $T \subseteq Q$ に対する処理は、まず、NIX部を用いてOID集合が検索される。次に、それらに対応するシグネチャファイルと問合せシグネチャの間に問合せ条件が成立するかを調べる。条件が成立すると判定されたターゲットオブジェクトについて、オブジェクトを検索して実際に問合せ条件を満たすか調べる。

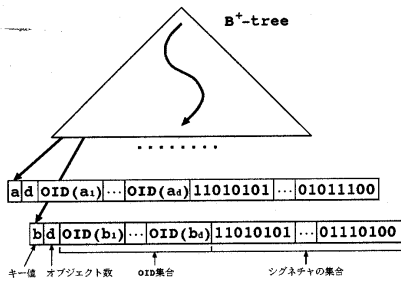


図2: S-NIXの構成

4 コスト解析

S-NIX, BSSF, NIXの3種類の索引機構について、コスト解析を行った。但し、 $T \supseteq Q$ における索引、 $T \subset Q$ におけるBSSFによる索引については、スマート検索方式[1]を評価した。

変数名	意味
V	集合要素の定義域の要素数
N	オブジェクトの総数
D_t	ターゲット集合の要素数
F	シグネチャのビット長
m	要素シグネチャにおける“1”の数(ウェイト)

表1: 記号とその意味

4.1 検索コスト

4.1.1 $T \supseteq Q$ の問合せの検索コスト

$V=16000$, $N=160000$, $D_t=10$ における検索コストを図3に示す。この場合、S-NIXのコストはNIXとほぼ同等で、BSSFに比べ優れたものになる。

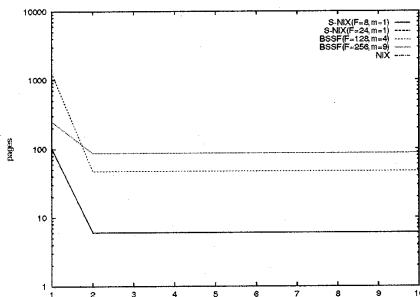


図3: $T \supseteq Q$ の検索コスト ($V=16000, N=160000, D_t=10$)

4.1.2 $T \subset Q$ の問合せの検索コスト

$D_t=10$ における検索コストを図4, 5に示す。S-NIXでは、 V, N が大きくなるほど、 B^+ 木部の絞り込みにより、BSSFに比べコストの増大が抑えられる。

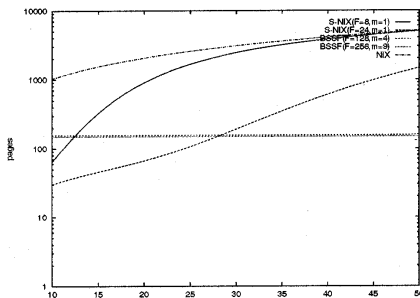


図4: $T \subset Q$ の検索コスト ($V=16000, N=160000$)

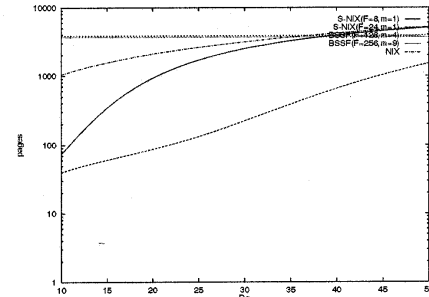


図5: $T \subset Q$ の検索コスト ($V=400000, N=4000000$)

4.2 格納コスト

$V=400000$, $N=4000000$, $D_t=10$ における格納コストを表2に示す。S-NIXの格納コストは、BSSF, NIXに対して劣っている。しかし、 F の長さを必要最小限にすることによって、コストの増大を抑えることができる。

D_t	ファイル	ページ数
10	S-NIX ($F=8$)	100463
	S-NIX ($F=24$)	133950
	BSSF ($F=128$)	23557
	BSSF ($F=256$)	39301
	NIX	80370

表2: 格納コスト ($V=400000, N=4000000$)

4.3 更新コスト

$V=400000$, $N=4000000$, $D_t=10$ における更新コストは、表3のようになる。S-NIXは、 F のサイズによりリーフノードサイズが大きくなり過ぎない限り、更新コストはNIXとほぼ同等に抑えられる。

D_t	ファイル	挿入コスト	削除コスト
10	S-NIX ($F=8$)	40	40
	S-NIX ($F=24$)	40	40
	BSSF ($F=128, m=4$)	36	3941
	BSSF ($F=256, m=9$)	77	3983
	NIX	40	40

表3: 更新コスト ($V=400000, N=4000000$)

5 おわりに

本研究では、シグネチャ付きNIXを提案し、そのコスト解析を行った。S-NIXは、 V, N の両方が大きくなるほど、シグネチャと B^+ 木を組み合わせた効果が発揮され、BSSF, NIXに対して検索コストの低減が実現された。

今後の課題としては、シグネチャの格納に関する冗長性の解決、圧縮を組み合わせたより効率的な索引構成法の検討等がある。

参考文献

- [1] Y. Ishikawa, H. Kitagawa, N. Ohbo. Evaluation of Signature Files as Set Access Facilities in OODBs, *Proc. ACM SIGMOD Conf.*, pp.247-256. May.1993.
- [2] 石川佳治, 北川博之, 大保信夫. シグネチャファイルによる集合値検索のコスト評価, *情報処理学会論文誌*, 第36巻第2号. Feb.1995.