

管理者に対しても機密を保持できる暗号化データベースの索引構成法

三浦 志保[†] 渡辺知恵美^{††}

[†] お茶の水女子大学理学部情報科学科 〒112-8610 東京都文京区大塚 2-1-1

E-mail: [†]shiho@db.is.ocha.ac.jp, ^{††}chiemi@is.ocha.ac.jp

あらまし 近年，データベース製品の管理運用を外部のデータベース技術者に委託するデータベースアウトソーシングサービスが普及しつつある．このときサービスの利用者は，機密事項であるデータを外部のデータベースサービスプロバイダに委託しなければならない．利用者はデータの内容をデータベース管理者に閲覧されたくないという要求を持つ．そのため，管理者に対しても機密を保持できる暗号化データベースシステムの研究が現在進められている [4]．このシステムでは，クライアント側でタプルごとにデータを暗号化し，それに索引を付与してサーバに置く．索引には，ドメイン領域を分割してそれぞれに割り当てられたバケット番号が格納される．検索時は以下の 2 段階問合せが行われる．まずフィルタリング処理として，タプルに付与された索引を用いてサーバで大まかな検索を行う．そして，該当したタプルに対してのみクライアント側で解精製を行う．文献 [4] では，バケットのデータ分布から索引情報が推測される可能性がある．よって本稿では，MaxDiff ヒストグラムに基づいた，より安全な索引構成法を提案する．
キーワード データベースセキュリティ，プライバシー保護，暗号，ヒストグラム

A Indexing Method for Encrypted Databases Protecting Confidential Data even from DBAs

Shiho MIURA[†] and Chiemi WATANABE^{††}

[†] Department of Information Sciences, Faculty of Science, Ocyanomizu University

E-mail: [†]shiho@db.is.ocha.ac.jp, ^{††}chiemi@is.ocha.ac.jp

Abstract Recently, the database out-sourcing service model which entrusts database management to database service provider has been gaining much attention. To use this model, user must entrust confidential data to external database administrator. Then users of the system require to protect their data from even database administrator. To resolve the problem, encrypting database system has been proposed [4]. Each tuple is encrypted and added privacy-preserving indexes. The indexes are created from bucket-id which divides their domain area. The query operator processes in the following two steps; (1) filtering query using privacy-preserving indexes in server, and (2) refinement query from decrypted tuples in client. The index method used in paper [4] is about the possibly of deduction from the bucket distribution. We propose a safety-higher bucketizing method based on MaxDiff histogram.

Key words database-security, privacy-preservation, cryptography, histogram

1. はじめに

近年，個人情報保護法の制定や大規模な情報漏洩事件が次々に明るみに出たことにより，データベースセキュリティに対する意識が急速に高まっている．多くの DBMS 製品は豊富なセキュリティ機能を備えているが，原則としてデータベース管理者に全面的な信頼を置くことを前提に提供されている．そのため，管理不十分による情報漏洩や管理者自身の内部犯行に十分に対処できない．

また，データベース製品の管理運用を外部のデータベース技

術者に委託するデータベースアウトソーシングサービスも現在普及しつつある．このときデータサービスの利用者は，委託業者の管理者に機密事項であるデータの内容を閲覧される可能性をなくしたいという要求を持つ．

そこで近年では，暗号化したデータベースと索引をサーバに置くことによって利用者側が管理者に頼ることなく情報漏洩を防ぐことのできる暗号化データベースシステムの研究が進められている [4]．暗号化データベースシステムでは，クライアント側でタプルごとにデータを暗号化し，セキュアな索引を付与し，それらをサーバに置く．検索時には，まずサーバで索引を

用いて大まかな検索を行ってから該当タプルをクライアント側で復号化し、解の精製を行う。このときタプルに付与する索引は、タプルの内容が侵入者に推測されず、かつ、サーバ側でのデータの絞込みが十分に行われるものでなければならない。これまで用いられている索引構成法では、属性のドメイン領域を等間隔のバケットに分割し各バケットにランダムに付けられたバケット番号をその属性の索引値とする方法が用いられている（その詳しい説明は2節に述べる。）しかしながらこの手法では属性値に偏りがある場合、各バケットに入るタプル数の関係から値が推測されてしまう可能性があり十分に安全であるとはいえない。

そこで本研究では、タプルの内容を侵入者に推測されにくいより安全な索引構成法を提案する。この索引構成法は、暗号化データベース内のテーブルを見ることができ、かつログも見ることもできるデータベース管理者に対しても機密を保持できるものでなければならない。本提案手法は MaxDiff を用い、タプルの挿入または変更時に、各バケットに入るタプル数が均等になるようにバケットを調整する。またダミーデータを効果的に用いることにより、ログや索引を分析しても値が推測できないようにする。本稿ではその初期的試みとして単一属性を対象にタプル挿入時のバケット調整法について提案しその検証を行った。まず第2節にて暗号化データベースについて説明した後、第3節で安全な索引構成法を提案し、バケット調整の素朴法について述べる。素朴法は各バケットに入るタプル数が均等にはなるものの実行時間などの問題がある。それらの問題についての検証を第4節で行った後、バケット内ヒストグラムを用いた解決法について述べ、実行時間が十分に短縮されることを検証する。第5節は関連研究、第6節はまとめと今後の課題である。

2. 暗号化データベースシステム

2.1 暗号化データベースに対する処理の流れ

図1にクライアントとサーバにおける、暗号化データベースに対する処理の流れを示す。データは全てクライアント側で暗号化されてから、サーバ側の暗号化データベースに格納される。これで例えばデータベース管理者であっても、タプルの内容を知ることはできなくなる（図1①）。暗号化データに対する検索は、データに付与したセキュアな索引を用いて後に説明する2段階問合せを行うことにより実現する（図1②）。またサーバ側には、暗号化データベースの他に索引情報（図1）を置き、クライアントがデータを挿入・更新したり検索したりする際には、自動的にここから情報を取得できるようにする。

2.2 暗号化と索引の付与

暗号化データベースのテーブルは、元データのタプルを暗号化した文字列と各属性の索引とで構成されている。リレーション R において、 A_i を属性、 A_i^S をその属性に対する索引、 $etuple$ を1つのタプルを暗号化した文字列とすると、下のよう

$$R(A_1, A_2, \dots, A_n) \quad R^S(etuple, A_1^S, A_2^S, \dots, A_n^S)$$

各属性の索引を作成するにあたって、まずは、その属性のドメイン領域を適当に分割する（以下、分割された1つ1つの領

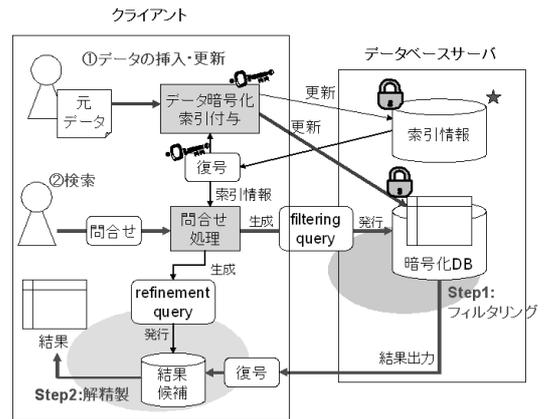


図1 暗号化データベースに対する処理の流れ

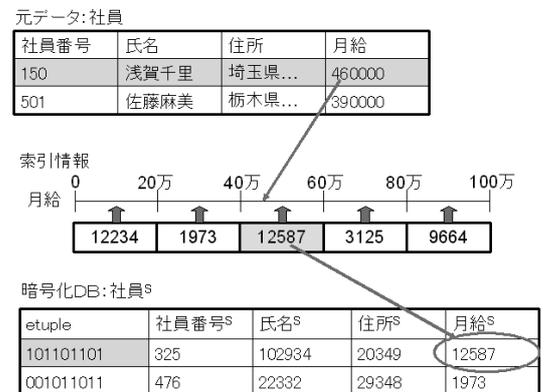


図2 暗号化と索引の付与

域のことをバケットという。）

$$\text{partition}(R.A_i) = \{p_1, p_2, \dots, p_k\}$$

次にそれぞれのバケット p_j にバケット番号をハッシュ関数により与え、その値 m を索引値として格納する。

$$\text{ident}_{R.A_i}(p_j) = m$$

また、それぞれのバケットの範囲、割り当てられたバケット番号は索引情報として別に格納する。索引情報もまたクライアント側で暗号化してから、サーバ上におく。

社員というテーブルで例を図2に示す。まず、社員（社員番号、氏名、住所、月給）社員^S($etuple$, 社員番号^S, 氏名^S, 住所^S, 月給^S)とし、社員^Sをサーバ側にある暗号化データベースに格納する。1行目の $etuple$ には“101101101”が格納されているが、これは“150 浅賀千里 埼玉県... 460000”を暗号化したものである。索引を作成する方法を月給を例に考えた場合、まず、 $\text{partition}(\text{社員.月給}) = \{(0, 20万], (20万, 40万], (40万, 60万], (60万, 80万], (80万, 100万)\}$ と領域を5つのバケットに分割する。そして、各バケットに $\text{ident}_{\text{社員.月給}}((40, 60万]) = 12587$ のようにバケット番号を与える。したがって、月給=460000は月給^S = 12587として格納される。

2.3 検索

ユーザが発行した問合せ文は、以下の2段階の問合せに分割され実行される。

- filtering query: サーバ側の索引に対する問合せ（フィルタリング処理）により、大まかな結果候補を抽出する。

- refinement query: クライアント上で復号化した結果候補 (フィルタリング処理により抽出) に対して解精製処理を行い, 最終的な結果を求める.

例えばクライアントが, 月給が 55 万円以上の社員の名前を求めるために「SELECT 名前 FROM 社員 WHERE 月給 \geq 550000」という問合せ文を発行したとする. その場合, まずはサーバから索引情報 (バケット情報) を自動取得してきて復号化を行い, その情報を元に filtering query 「SELECT etuple FROM 社員^S WHERE (月給^S = 12587 or 月給^S = 3125 or 月給^S = 9664)」を生成する. この filtering query をサーバに発行することによって, 月給が 55 万円以上の可能性のある社員のタプルを検索する. 上で述べた通り, この結果は大まかな結果候補にすぎない. filtering query では「WHERE (月給^S = 12587 or 月給^S = 3125 or 月給^S = 9664)」としているので, 厳密には月給が 40 万円以上の社員のタプルが返ってくるからである. 本当の結果を求めるために, このフィルタリングされた結果候補をクライアント側で復号化した後に, refinement query 「SELECT 名前 FROM 社員 WHERE 月給 \geq 550000」を発行する. この作業によって, 真の結果を得ることができる.

この 2 段階の問合せの仕組みの大きな特徴は, サーバ側でフィルタリング処理をすることによりクライアント側での処理が少なくてすむため, 検索効率をさほど下げることなく安全性を保てるという点である.

3. 安全な索引構成法の提案

3.1 従来の索引構成法による安全性

暗号化タプルに付与する索引は, タプルの内容を侵入者が推測できないものでなければならない. しかし, 文献 [4] で用いられている等間隔の領域分割法では, 元データにおいて同じ値 (または近い値) 同士のもは同じバケットに入るため, 暗号化データベース内でも同じ値で格納される. さらに, それぞれのバケットに入っているタプル数は一目で分かる. したがって, それぞれのバケットに入っているタプル数に偏りがある場合, 分布解析などによりデータの内容を予測されてしまう可能性がある. 例えば給与の場合, 高額な給与をもらっている人はそれほど多くないので, 各バケットに入っているタプル数には偏りが生じる.

サーバ側にたくさんのサンプルが与えられたとしても, 予測できない索引構成方法を考える必要がある. そこで我々は MaxDiff [7] による領域分割をベースにした索引構成法を提案する.

3.2 索引構成法の提案

3.2.1 基本概念

本研究では, MaxDiff [7] での領域分割法を用いる. これは, 1 つのバケットに入るタプル数の差が閾値以下になるように分割する方法である. この方法を用いることにより, バケットによるタプル数の偏りが生じなくなるため, 分布解析からデータの内容を予測されるのを防ぐことができる.

具体的には, データが挿入される度にそれぞれのバケットの範囲を変えることにより, それぞれのバケットに入るタプル

数に一定数以上の差を生じさせないようにする. 以下, それぞれのバケットに入るタプル数の差の上限を MaxDiff と表す. バケットを分割する際, バケットの個数の上限を決めておく方法と, バケットに入るタプル数の上限を決めておく方法が考えられるが, 本稿ではバケット分割にかかるコストを考えバケットの個数の上限を決めておく方法について述べる. 以下, バケットの個数の上限を p_{max} と表す. 各バケットに入るタプル数の上限を決めた分割法の提案や, 本提案法との比較は今後の課題として行う予定である.

3.2.2 データ挿入時のバケット調整法

データ更新の際に使う, 索引情報テーブルの例を表 1 に示す. ただし, 表 1 は図 2 に示したテーブル '社員' の索引情報から, 属性 '月給' の部分のみを示している. 下限値, 上限値とはあるバケットにおける上限と下限を表しており, バケット番号とはそのバケットに割り当てられている値である. また, 表 1 に示しているのはクライアント側で復号化した後の状態であり, 実際にサーバ側に格納しているときは, 下限値・上限値・バケット番号・タプル数はまとめて暗号化した文字列となっている.

テーブル名	属性	下限値	上限値	バケット番号	タプル数
社員	月給	0	20 万	12234	19
社員	月給	20 万	40 万	1973	20
社員	月給	40 万	60 万	12587	21
社員	月給	60 万	80 万	3125	21
社員	月給	80 万	100 万	9664	19

表 1 索引情報テーブル

クライアントがデータ挿入したときの処理の流れを以下に示す.

- (1) サーバから索引情報を取得し復号化する.
- (2) データ挿入後, 必要ならばバケット調整を行い, 索引情報および暗号化データベースを更新する.

上記の (2) における, 具体的な更新方法について説明する.

タプル r をリレーション R に挿入する際, リレーション R の各属性 $R.A_i (1 \leq i \leq n)$ に対して, $\text{partition}(R.A_i) = (p_0, p_1, \dots, p_k)$ の領域を以下の手法に従って調整し, 属性 $R.A_i$ の値がバケット $p_j (1 \leq j \leq k)$ の範囲に含まれるタプルの数 $\text{num}(p_j)$ の差が一定数以下となるようにする. 以降 $R.A_i$ の領域調整について述べるため, 「属性 $R.A_i$ の値がバケット p_j の範囲に含まれるタプル」を, 単に「 p_j の該当タプル」と呼ぶことにする. 各バケット p_j は (下限値, 上限値] と表し, まだデータが入っていない状態でも, 初期値として, $p_0 = (R.A_i \text{ のドメインの下限値}, R.A_i \text{ のドメインの上限値}]$ が与えられる. $\text{bottom}(p_j)$ はバケット p_j の下限値, $\text{top}(p_j)$ はバケット p_j の上限値を示すものとする.

今, MaxDiff=2, $p_{max}=5$ として, “insert into 社員 values ('200', '柳平有美', '東京都...', '610000')” が発行されたときの更新方法を属性 '給与' を例に示す.

- (1) バケットの個数 $< p_{max}$ のとき
分割数がまだバケット数の上限に達してないときは, 基本

的に1バケットに1タプルとなるようにデータ挿入のたびにバケットを分割する。上記の例では属性値‘61000’が該当するバケットを探し、図3のように、そのバケット p_3 を $(\text{bottom}(p_3), 61000]$ と $(61000, \text{top}(p_3)]$ に分割する。2分割したうち小さい方のバケットに‘61000’を格納する。結果各バケットの上限値のところのみデータが存在する状態となる。また挿入する値がバケットの上限値である場合は分割できないためそのまま挿入する。

挿入後、バケット内タプル数の最大値とバケット内タプル数の最小値との差が MaxDiff を超える場合（例では、各バケットの該当タプル数に3以上の差が出てしまった場合）がある。そのときは該当タプル数が最小であるバケットにダミータプルを挿入することで調整を行う。

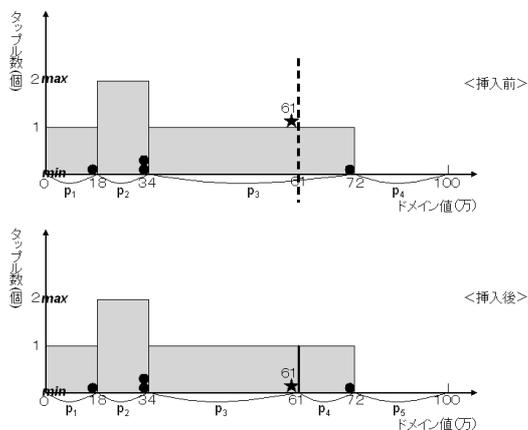


図3 バケットの2分割

(2) バケットの個数 = p_{max} のとき

バケット数がその上限 p_{max} を満たしている場合も、まずは挿入する値‘61000’が該当するバケットを探し、そのバケットに格納する。格納後、バケット内タプル数の最大値と最小値との差が MaxDiff を超えない場合は、バケット調整は行わない。

タプルの挿入後に、バケット内タプル数の最大値と最小値との差が MaxDiff を超えてしまう場合、該当タプル数が最大であるバケット p_j に対しバケットの切り崩しを行う。

[バケットの切り崩し]

バケット p_j の隣接バケットを見て、該当タプル数の少ない方のバケット p_i ($p_i = p_{j-1}$ または p_{j+1} とする) へ p_j 内のタプルを切り崩す。該当タプルの切り崩しは p_j および p_i の範囲を変えることによって行う。 $p_i = p_{j+1}$ である場合はバケット p_j をスキャンし、 p_j 内で2番目に大きな値を p_j の上限値および p_{j+1} の下限値とすることによって p_j 内の最大値を p_{j+1} に移動する。 $p_i = p_{j-1}$ である場合はバケット p_j 内の最小値を求め、 p_{j-1} の上限値および p_j の下限値とすることによって p_j 内の最小値をバケット p_{j-1} に移動する。

ただし、 p_j の範囲を変更することにより、 p_j の該当タプル数が最小になってしまう場合はその処理を行わず、逆隣への切り崩しを同様に試みる（これを以下、方向転換と呼ぶ）

p_j の切り崩しにより、隣接バケット p_i の該当タプル数が最大

となる場合は、 p_i に対して切り崩し処理を行う。

この切り崩し処理を、各バケットの該当タプル数の差が MaxDiff 以下となるか、切り崩し処理がそれ以上できなくなるまで繰り返す。

図4ではバケット p_3 からバケット p_2 への切り崩しをするために、 p_3 内の最小値である53を p_2 の上限値および p_3 の下限値とする。しかし、この処理により p_2 の該当タプル数が4となり、 p_5 の該当タプル数との差が $\text{MaxDiff}=2$ を超えるため、 p_2 に対して切り崩し処理を行い、 p_1, p_2 の範囲を変えることによって、 p_2 の該当タプルの内1つを p_1 に渡している。

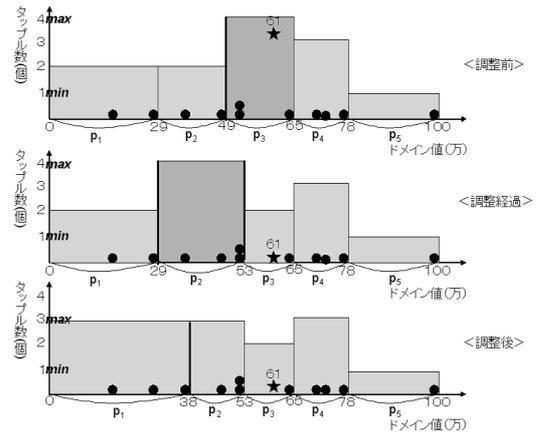


図4 バケットの切り崩し

図5は、方向転換が必要な例である。まずバケット p_4 からバケット p_3 への切り崩しを試みるが、 p_4 内の最小値を持つタプル数は3個あり、これらを p_3 に渡した場合、 p_4 の該当タプル数は明らかに最小になってしまう。したがって p_4 から p_3 への切り崩しをせずに方向転換をし、バケット p_5 の方へ p_4 内の最大値を持つタプルを渡している。

該当タプル数最大のバケットの切り崩しを上記の通りに行った後も、その差が MaxDiff 以下にならない場合が考えられる。その場合、該当タプル数が最小である領域 p_j を見つけ、バケットの埋め合わせを行う。

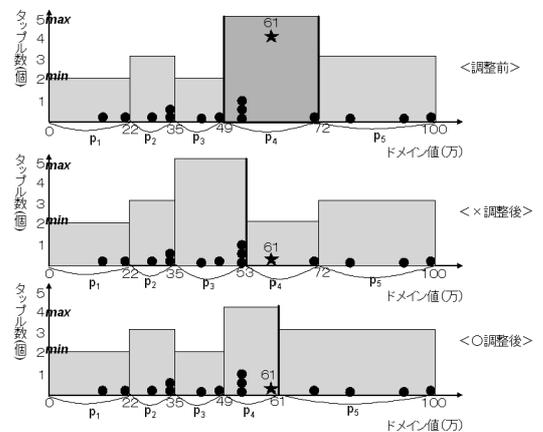


図5 方向転換を含むバケットの切り崩し

[バケットの埋め合わせ]

該当タプル数最小のバケット p_j の隣接バケットのうち該当

タプル数の多い方からバケット p_j へタプルを切り崩すことにより、 p_j の該当タプル数を増やす。ただし、切り崩しの時と同様、タプルの受け渡しにより隣接バケットの該当タプル数が最小となる場合は方向転換を行う。

[ダミータプルの挿入]

バケットの切り崩しと埋め合わせを行っても各バケットの該当タプル数の差が MaxDiff を超える場合は該当タプル数が最小であるバケットにダミータプルを挿入することで調整をはかる。

本システムでは、ダミータプルが増えないための工夫が施されている。1つ目は、ダミータプルの上書きである。すでにダミータプルが格納されているバケットに新しい値が格納された場合、そのダミータプルを上書きする。そうすることで無駄にデータ量を大きくすることを防いでいる。2つ目は、ダミータプル数をそのバケットの該当タプル数として数えないことである。バケット内タプル数の最大値と最小値との差が MaxDiff を超える場合、切り崩し処理や埋め合わせ処理によりバケット調整を行うが、始めに該当タプル数を数える時点で、ダミータプルを数えないことで無駄なバケット調整を行うのを防ぐことが出来る。よってダミータプルが必要以上に増えるのを防ぐことが出来る。

4. 提案手法の改良

4.1 提案手法における懸案事項

前章で提案した索引構成法には大きく分けて4つの問題点が考えられる。

(1) バケット p_j の該当タプル数を調整する場合、 p_j 内の最小値もしくは2番目に大きな値を求めるために、 p_j 内のタプルデータを全てサーバ上の暗号化テーブルから取得し復号化しなければならない。よって、サーバから取得するタプル数が膨大だった場合、実行時間も膨大となる可能性がある。

(2) 1個のデータを挿入する度に、バケット調整においてバケット番号が変わった複数のタプルを更新しなければならない。更新タプル数が膨大になった場合、実行時間も膨大となる可能性がある。

(3) 本来、各バケットに与えられるバケット番号はランダムなため、サーバ側ではバケットの隣接関係は分からないようになっている。しかしデータ挿入に伴うタプルの更新パターンによって、バケットの隣接関係が推測できる可能性がある。

(4) 各バケットにおける該当タプル数調整のために、ダミーデータが増えて必要以上にデータ量が大きくなる可能性がある。

4.2 検証

前節で挙げた4つの問題点の影響を調べるため、それぞれ4つの実験を行った。4つの実験において挿入したデータは、いずれも一様乱数により生成した値である。また、属性のドメインは $[0,50000]$ 、最大バケット数 p_{max} は20、 maxdiff は2である。

実験(1) サーバから取得するタプル数と時間の連動性

1000回のデータ挿入を行い、各挿入処理時バケット調整のためにサーバから取得したタプル数と実行時間を計測

する。

実験(2) 更新タプル数と実行時間の連動性

1000回のデータ挿入を行い、データ挿入に伴うバケット調整による更新したタプル数と実行時間を計測する。

実験(3) タプルの更新パターン

1000回のデータ挿入を行い、データ挿入に伴う各バケットの更新タプル数を調べる。

実験(4) ダミーデータ

100回のデータ挿入を行い、各バケットのダミータプル合計数を調べる。これを10回繰り返して、ダミータプル合計数の平均を求める。

図6に実験(1)の結果を示す。図6(a)が実行時間、図6(b)がサーバから取得したタプル数のグラフである。グラフからわかるように、サーバから取得したタプル数が多いとき、実行時間も大きくなっている。これは暗号化テーブルを全てスキャンし、該当するタプルを復号化するのにコストがかかっているからだと考えられる。

図6(c)に実験(2)の結果(更新タプル数)を示す。更新タプル数の増え方と実行時間の増え方にはずれが生じている。このことから、実行時間が更新タプル数よりサーバから取得するタプル数に依存していることが分かる。サーバから取得しなければならないタプル数は普通、更新タプル数よりも大幅に多い上に、復号化作業に伴うからである。

図7に実験(3)の結果の一部を示す。1000回のデータ挿入のうち37回目から46回目のデータ挿入の結果を示しており、またバケットも全20個のうち p_2 から p_{12} までの11個のみ示している。ただし、示されていないバケットに関してはタプルの更新がなかったものとする。図から37, 38, 39, 41, 43回目の挿入では、タプルの更新がないことが分かる。これはデータの挿入後、バケット調整を行う必要がなかった場合だと考えられる。また40, 42, 44回目の挿入では、データを挿入したバケットの隣のバケットのタプルが更新されている。これは挿入後、1回バケット調整が行われた場合だと考えられる。また45, 46回目の挿入では、データを挿入したバケットの隣のバケットから遠いバケットまで連続してタプルが更新されている。これは挿入後、複数回バケット調整が行われた場合だと考えられる。

全体(1000回)を通して最も多いのは、データを挿入したバケットの隣のバケットのタプルが更新されるパターンである。またサーバにはどのバケットにデータが挿入されたかがログとして残ってしまう。そのためタプルが更新されたバケットと、ログから分かるデータを挿入したバケットとが隣接関係にある可能性が高いと予測することができてしまう。

図8に実験(4)の結果を示す。図8(a)は10回分の施行におけるダミータプル数の合計値である。trial3, 6, 30以外はほとんどダミータプルがなく、均等に分けられていることがわかる。しかし trial3, 6, 30においてはダミー数が非常に多くなった。最もダミータプルの多い trial3 における各バケットのダミータプル数を図8(b)に示す。この場合は p_1 から p_7 のタプル数に合わせる形で p_8 以降のバケットにダミータプルを含めている。本手法ではタプルを挿入するたびに差分的

にバケット調整を行っているが、挿入する値に偏りがある場合に切り崩しがうまくいかずダミータプルが増える結果になると思われる。また、バケット数 p_{max} を小さく設定したり、許容するバケットごとの該当タプル数の差 MaxDiff を小さく設定したりした場合は、バケット調整による高さの均等化が難しくなり、ダミータプル数が大幅に増える可能性がある。また、分布的に偏りのあるデータを挿入していった場合も同様にしてダミータプル数が増える可能性がある。

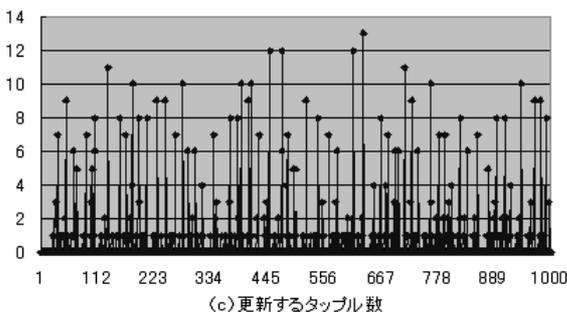
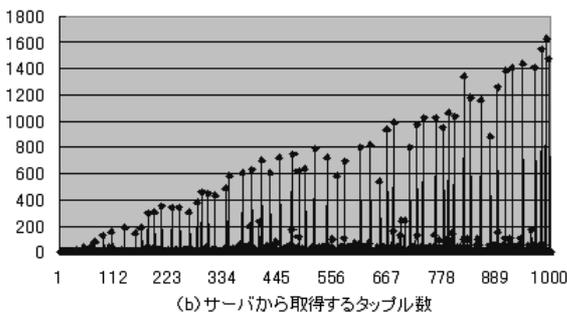
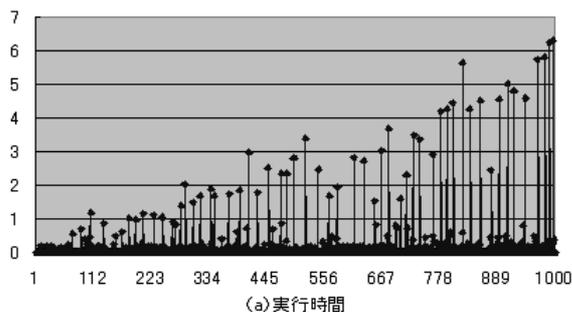


図 6 素朴法の実験結果 (a) 実行時間, (b) サーバから取得するタプル数, (c) 更新するタプル数

挿入回数	挿入バケット	P_2	P_3	P_4	P_5	P_6	P_7	P_8	P_9	P_{10}	P_{11}	P_{12}
37回目	P_9	0	0	0	0	0	0	0	0	0	0	0
38回目	P_6	0	0	0	0	0	0	0	0	0	0	0
39回目	P_{10}	0	0	0	0	0	0	0	0	0	0	0
40回目	P_{10}	0	0	0	0	0	0	0	0	0	1	0
41回目	P_{12}	0	0	0	0	0	0	0	0	0	0	0
42回目	P_6	0	0	0	1	0	0	0	0	0	0	0
43回目	P_7	0	0	0	0	0	0	0	0	0	0	0
44回目	P_7	0	0	0	0	0	0	1	0	0	0	0
45回目	P_{10}	0	1	1	1	1	1	1	1	0	0	0
46回目	P_9	1	1	1	1	1	1	1	0	0	0	0

図 7 タプルの更新パターン

試行	ダミー数
trial1	7
trial2	7
trial3	43
trial4	8
trial5	4
trial6	17
trial7	2
trial8	0
trial9	0
trial10	30
sum	11.8

(a)ダミー数合計

	下限値	上限値	タプル数	ダミー数	ダミー数を除いたタプル数
p[1]	0	8648	8	0	8
p[2]	8648	10304	8	0	8
p[3]	10304	13693	8	0	8
p[4]	13693	17010	7	0	7
p[5]	17010	20750	7	0	7
p[6]	20750	34779	7	0	7
p[7]	34779	27481	7	0	7
p[8]	27481	30948	7	3	4
p[9]	30948	33360	7	1	6
p[10]	33360	36300	7	1	6
p[11]	36300	37552	7	3	4
p[12]	37552	38763	7	3	4
p[13]	38763	42640	7	1	6
p[14]	42640	43184	7	5	2
p[15]	43184	46012	7	3	4
p[16]	46012	46579	7	5	2
p[17]	46579	48591	7	1	6
p[18]	48591	49556	7	6	1
p[19]	49556	49627	7	6	1
p[20]	49627	50000	7	5	2
				43	

(b) trial 31における各バケットのダミー数分布

図 8 ダミーデータ

4.3 解決方法

前節に述べた各問題点に対する解決方法の方針を以下に示す。バケット調整にかかる時間の短縮 各バケット内のヒストグラムを索引情報に格納し、それを用いてバケット調整を行うことによりサーバからのデータ取得および復号化にかかる処理を減らす。

隣接関係が推測されないバケット更新 バケット更新時に、更新しないバケットのうち該当タプル数が最大でないバケットにダミータプルを挿入する。

ダミータプル数の削減 前節に述べたようにダミータプルが増えてしまうのはタプルを挿入するたびに差分的にバケットの調整を行っていることが一つの原因として考えられる。削減法としては、バケット数や該当タプル数の差 MaxDiff の動的な調整、また、定期的にバケット分割の組み直しなどを行うことによってダミータプル数を削減する方法が考えられる。

これらの解決法のうちヒストグラムを導入した改良法について次節に説明しその検証を行う（他の解決法については今後の課題とする。）

4.3.1 ヒストグラムを導入した索引構成法

実験結果より、実行時間にはサーバから取得するデータ量が非常に影響することが分かった。よって、サーバから取得するデータの量を最小限に抑えるため、3章で提案した索引構成法にさらにヒストグラムを導入する。

[ヒストグラムの生成]

図 9 に示すように、それぞれのバケット内を等間隔の bin に分割する。bin の幅はあらかじめ決めておくこととする。それぞれの bin に該当するタプル数を文字列にし、属性”histogram”を索引情報に付加し文字列を格納する。図 9 の例では、属性”histogram”の値は”0:3,1:7,3:2,4:7,5:5”というように左から順にヒストグラム番号を付け、ヒストグラム番号と該当タプル数のペアを並べる。

[ヒストグラムを用いたバケット調整]

正の方向に切り崩しを行う場合、切り崩しを行うバケット p_i 内で 2 番目に大きい（右から 2 番目の）bin の最大値をバケット p_i と p_{i+1} の新しい境界とする。また負の方向に切り崩しを行う場合、切り崩しを行うバケット p_i 内で最小の（1 番左の）

bin の最大値をバケット p_i と p_{i-1} の新しい境界とする．図 10 は負の方向への切り崩しの例を示したものである．

このように，ヒストグラムを利用してバケット調整を行うことにより，新しい境界線を決定するためにサーバからデータを取得する必要がなくなる．ヒストグラムの情報は索引情報に格納されているので，索引情報を復号化し見るだけでバケット調整ができる．また 1 個のデータ挿入の際に行われるバケット調整は 1 回とは限らず，複数回ある場合がある．このとき 1 回の調整の度にサーバ上にある暗号化テーブルを更新するのでは無駄な時間がかかってしまう．調整の度に，ヒストグラム情報を含む索引情報だけを更新し，バケットの範囲変更の情報は蓄積していく．そして全ての調整が終了してから暗号化テーブルを更新する．そうすることで，更新タプル数も最小限に抑えることができる．

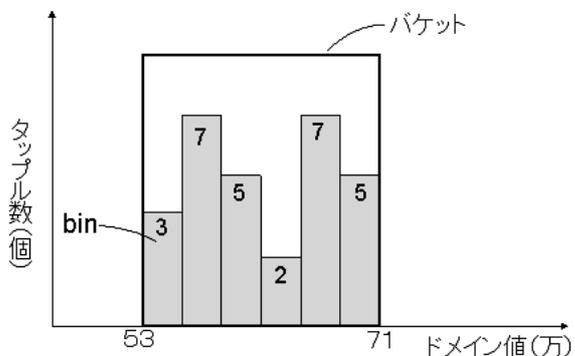


図 9 ヒストグラム生成

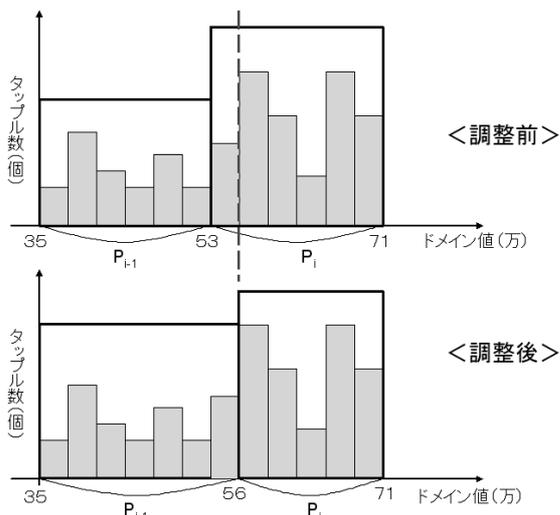


図 10 ヒストグラムを用いたバケット調整

4.3.2 改良法の検証

ヒストグラムを導入した索引構成法において，ヒストグラム導入前の実験と同じ条件でデータ挿入を行い，その実行時間を計測した．ただし，バケット内の bin の幅は 20 とした．図 11 はヒストグラム導入後の実行時間 (a)，取得タプル数 (b)，更新タプル数 (c) である．

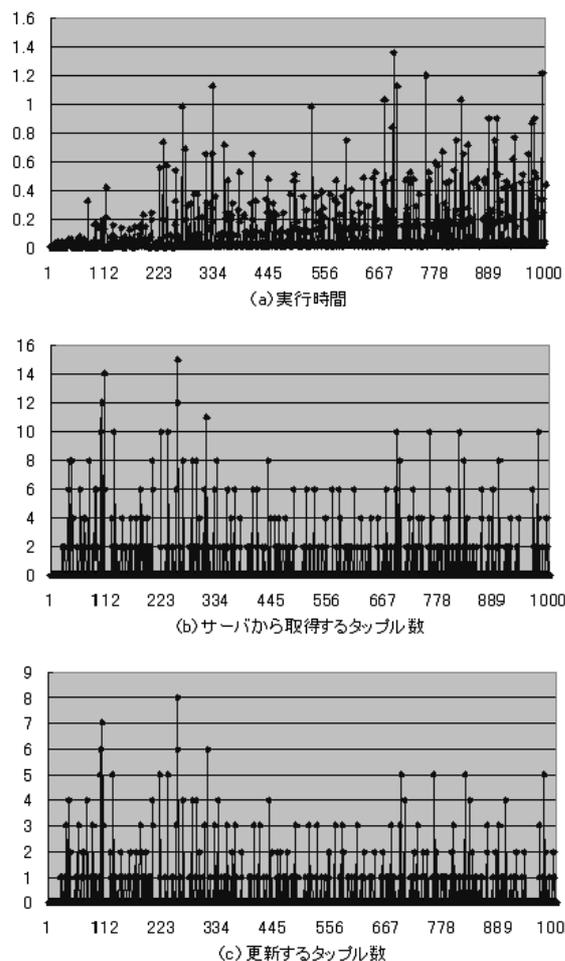


図 11 ヒストグラム導入後の実験結果 (a) 実行時間, (b) サーバから取得するタプル数, (c) 更新するタプル数

図6と比べてみると、ヒストグラム導入後はサーバからの取得タプル数と更新タプル数に大きな差がなくなっており、更新タプルに応じて必要なだけのタプルがサーバから取得されるようになっていたことが分かる。そのため、実行時間も大幅に短縮され、最大時間で比べると3分の1から4分の1程度にまで小さくなっている。このことから、ヒストグラムの導入は実行時間短縮に有用であると考えられる。

5. 関連研究

データベースにおけるプライバシー保護やセキュリティの研究は、従来よりデータベースの主要な研究分野として数多く行われてきた[3]。近年ではプライバシー漏洩の問題に対する推論制御、統計データベースなどが注目されている。

また、一部の情報を隠蔽したり暗号化することで十分にプライバシーを保護したまま問合せやデータマイニングをするための研究も盛んに行われている[8]。Agrawalらの提案したReconstruction-Based Approach[1]は個人情報に意図的にランダムノイズを乗せてプライバシーを保護し、ベイズの定理を用いて、逐次的に確率密度関数を推定し、決定木によるマイニングを実現する。また、k-anonymity[5]はデータを開示する際、一部の情報を抽象化したり隠蔽することで同じ情報を持つ人数を一定数以上に制御する手法である。

暗号化データベースも個人情報をサーバ管理者に対して保護しながら問合せを行う手法であり、暗号化したタプルに付随する索引を用いて2段階の検索を行い、問合せを実現している。本手法は文献[4]によって提案され、その後集約演算に対する処理やアクセス制御への適用が研究されている。安全な索引構成法に関しては文献[2]にてEqui-depthによる分割法を利用し、処理効率と安全性を考慮した最適なバケット数を求める提案がされている。しかしながらこれらの研究ではデータはあらかじめすべて用意されているという前提であり、更新が起こることを想定していない。複数のユーザーで暗号化データベースを共有する場合には更新が頻繁に行うこともあると考えられるため、我々の手法も必要となると考えられる。また適切なバケット数の導出法に関して文献[2]を参考に更新に対応できるように改良する必要があると考えられる。

6. まとめと今後の課題

本稿では新しい索引構成法の1つとして、MaxDiffでの領域分割法を提案した。本提案手法はバケットの個数を固定し、バケット調整により高さを均等化する方法をとっている。また、ヒストグラムを導入し1回のバケット調整の度にサーバ上の暗号化テーブルをスキャンする必要をなくし、計算量を小さくする方法を併せて提案し、その有用性を実験により確認した。

本稿で行った試みは安全な索引を構成するための初期的な試みであり、今後多くの課題が残っている。例えば、本提案手法は、年齢や月給等の特定の範囲のある数値データには適用できるが、住所、氏名のようなデータには効率よく適用できないため、範囲を持たないデータに対する対処法を考えなければならない。また、削除や更新に伴うバケット調整、適切なバケット数やバケット内タプル数の検証やそれらを動的に変更するアル

ゴリズムの検討も必要である。また、今回は1属性のみを扱ったがテーブル内にある複数の属性に対する方法も考えなければならない。また、今回の索引構成法は部分一致による検索に対処することができないことから、それらに対処する方法も考えていく必要がある。また、今回は提案手法の懸念事項の解決法としてヒストグラムを用いた更新時間短縮の検証のみを行ったが、ほかの解決法も試みる必要がある。また、本提案手法ではバケット調整による索引値の変化をデータ挿入ごとに観測した場合、バケットの隣接関係がわかってしまう。例えば、給与の索引値が12234であったタプルがテーブルの更新後に1973という索引値に変わった場合、バケット番号が12234のバケットと1973のバケットは隣接関係にあることがわかってしまう。これを解決するためには、

このように、どの更新がバケット調整により生じた真の更新かを分からなくすることで隣接関係が推測されないようにできると考えられるので、検討していきたい。

今後はこれらのシステムを実装、改良し、安全性やコストの面から比較・検証していきたい。

文 献

- [1] R. Agrawal and R. Srikant.: "Privacy-Preserving Data Mining," *Proceedings of ACM SIGMOD Conference on Management of Data*, pp.439-450, 2000
- [2] Hore B., Mehrotra S. and Tsudik G.: "A Privacy-Preserving Index for Range Query," *Proceedings of the 30th VLDB Conference*, pp.720-730, 2004.
- [3] Maria Grazia Fugini, Silvana Castano(Editor), Giancarlo Martella (Contributor): "Database Security", *ACM Press Books/Addison-Wesley Publishing*, p.456, 1995.
- [4] H. Hacigumus, B. Iyer, C. Li, and S. Mehrotra.: "Executing SQL over Encrypted Data in the Database-Service-Provider Model," In *Proceeding of the 2002 ACM SIGMOD International Conference on Management of Data*, pp. 216-227, June 2002.
- [5] Sweeney L.: "K-anonymity: A Model for Protecting Privacy," *International Journal on Uncertainty, Fuzziness and Knowledge-based System*, 10(5), pp.557-570, 2002.
- [6] G. Piatetsky-Shapiro and C. Connell.: "Accurate estimation of the number of tuples satisfying a condition," In *Proceeding of ACM SIGMOD*, 1984.
- [7] V. Poosala, P. J. Haas, Y. E. Ioannidis and E. J. Shekita.: "Improved histograms for selectivity estimation of range predicates," *Proceedings of ACM SIGMOD*, pp.294-305, 1996.
- [8] Verykios V. S., Bertino E. and Fovino N.: "State-of-the-art in Privacy Preserving Data Mining," *SIGMOD Record*, 33(1), 2004.
- [9] M. Winslett and J. D. Ullman.: "Jeffrey D. Ullman speaks out on the future of higher education, startups, database theory, and more," *SIGMOD Record*, 30(3), 2001.