

データベース用自然言語インタフェース

5Q-1

「dBmate」

伊藤 篤 高橋 清一

㈱CSK総合研究所

1. はじめに

データベースの使用は初心者にとっては困難であり、そのため、自然言語インタフェースを持ったDBMSがいくつか考案されている。しかし、自然言語から派生する諸々の問題により実用レベルのものはないというのが現状であろう。

そこで、我々は、パソコンベースで普及しているリレーショナル・DBMSを対象とし、実用に耐えるシステムを目指して、自然言語インタフェース「dBmate」を開発した。このシステムは、自然言語からDBMSのコマンドを生成・実行し、自然言語による情報検索、情報操作を可能にしている。

2. システム概要

dBmateでは、自然言語によって生じるデメリットを最小限にして、自然言語のメリットを生かし、実用レベルにもっていくために、以下の点を工夫した。

- ・自然言語システムにありがちな問題点である「まず最初に目的に合わせて辞書を作成しなければならない」という面倒をなくす。これは、後述のDB依存辞書の自動作成機能によって達成した。
- ・C言語によるパーザを使い解析を高速にする。
- ・分かち書きが要らない。
- ・漢字かな混じり文でユーザの直感に近い文を入力できる。
- ・未定義語の処理によって、全ての使用する単語を辞書に登録する必要はない。
- ・「営業の男は？」の例のようにかなり省略した文も解

析・実行する。

- ・解析が失敗すると間違ったと思われる場所に↑を表示して再入力を促す。
- ・履歴機能により過去に入力した文はすぐに呼び出してエディット・実行ができる。
- ・例文メニューによりよく使う文は簡単に入力できる。
- ・データベースを自動的に選択・結合する。

1) システム構成

システム構成を図1に示す。dBmateは3層の構造を持つ辞書(文法)とデータベースの状態を保持する管理テーブルを持つ。

ユーザがかな漢字変換FEPによって入力した自然言語文から、辞書とパーザの構文解析によって解析木を作成する。さらに解析木は意味解析によってSTD(Semantic Tree for DBMS)と呼ぶ中間表現に変換される。ユーザの入力した自然言語文は省略が多かったり、データベースの指定がなかったりするので、このままではSTDは完全ではない。そこで、足りない情報を補ってSTDを完全にする。この情報は主にフィールドとファイル(データベース名)である。

STDが完全になるとコマンド生成部によってDBMSプログラムに変換しDBMSに引き渡して自動的に実行される。結果はDBMSによって画面・プリンタ等に出力される。

3. データ構造

1) 辞書の構造

辞書はシステム辞書、DB依存辞書、ユーザ辞書からなる。

システム辞書はシステムに組み込まれている辞書で、文法に関する情報と必要最低限の単語(助詞等)を含む。これは固定であり一切変更されない。

DB依存辞書はシステム立ち上げ時に指定したディレクトリに存在する(複数の)データベースから自動的に作成する。ファイル名、フィールド名、よく使われそうなデータ等を登録する。DB依存辞書の自動作成と同時に管理テーブルも作成する。

システム辞書とDB依存辞書によって一通りの操作が可能であるが、さらにユーザ独特の表現を可能にするためにユーザ辞書がある。ユーザ辞書はユーザが必要に応じて付け足すもので「XXXを辞書に登録する」といった文を入力することやメニューによって登録する。

2) STD

パーザによって作られた解析木から意味解析によってSTDを作成する。STDの構造を図2に示す。STDのルートノードは[FRAME]と呼ばれ6つの要素からなる。最初の<機能名>にはCREATE/USE/APPEND/EDIT/DELETE/DISPLAY/INDEX/REPLACE/AVERAGE/COUNT/SUM/RECALL/QUIT等が入る。これはDBMSのコマンドに1対1に対応するのではなく、データベースを作りたい(CREATE)とか情報を得たい(DISPLAY)といった機能に対応する。実際のどのコマンドを生成するかはコマンド生成部の仕事である。

<機能名>以外は子ノードとして[SUBFRAME]を持つ。[SUBFRAME]は<operation><left><right>からなるセルで

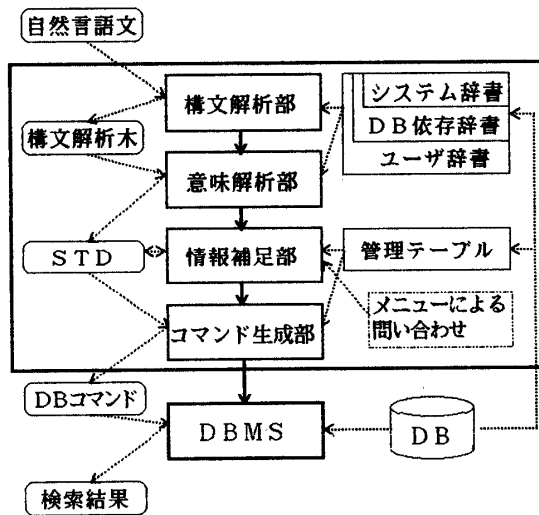


図1 システム構成

「dBmate」

Natural Language Interface for DBMS  
Atsushi ITOH, Seiichi TAKAHASHI  
CSK Research Institute

<left><right>にはさらに[SUBFRAME]が入る。<operation>に 0 が入る場合は特別にリーフセルを意味し、このときは<left>にそのセルのタイプが入り<right>にはデータ等が入る。

例えば「営業の男は?」という文が入力された場合のSTDを図3(a)に示す。

3) 管理テーブル

管理テーブルにはデータベースの状態が記述される。データベースの状態とは、データベース名、そのデータベースの作業領域番号とそのデータベースの持つフィールドの状態である。フィールドの状態はフィールド名、フィールドの型(文字型/数値型/論理型/日付型)、インデックスの状態からなる。インデックスの状態はさらにインデックスファイルを管理する番号、状態(有り/無し/不整合)、マスタキーかどうか、データの頻度(異なったデータが多い/普通/同じデータが多い)からなる。システムはこれらの情報からSTDの足りない部分を補足する。

4. 意味解析後の処理概要

1) フィールド情報の補足

STDの中でフィールド指定のないデータを見つけると、フィールド情報を補う。管理テーブルからフィールド名を抽出してユーザにメニュー形式で問い合わせる。ユーザが選択するとフィールドが決まる。

2) ファイル情報の補足

フィールド情報が充足すると、次はどのデータベースから検索するかを決める。フィールドの指定によっては2つのデータベースを結合して情報を得ることも有り得る。場合によってはファイルの選択が一意に決まらない場合がある。このようなときはユーザにメニュー形式で問い合わせて確定する。

フィールド情報、ファイル情報を補足したSTDの例を図3(b)に示す。

3) コマンド生成

STDが完全になると、<機能名>によって各コマンド生成を行なう。コマンドが生成されるとDBMSに引き渡され実行される。ここでコマンドといっても一つのコマンドではなくコマンド列のことである。プログラムといっても良い。「営業の男は?」の場合に生成されるプ

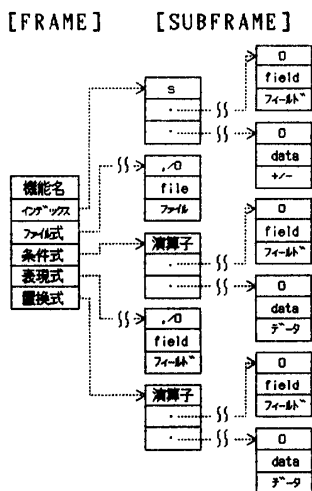


図2 STDの構造

ログラムは図3(c)を参照されたい。実行が終了すると再び自然言語文の入力待ちになる。

5. おわりに

d B m a t e はユーザが使いやすいシステムとなるために、DBMSに合わせてd B m a t e 側でいろいろと工夫した。しかし、真に使いやすいシステムを考えるならば、DBMSの側から歩みよるべきではないだろうか? 例えば、ファイル(データベース)の自動選択・連結等は自然言語から派生した要求とはいえ、DBMS側でやるべきことであると思う。

我々は単なるプロトタイプではなく、実用的なシステムを目指してd B m a t e を開発した。その目的はほぼ達したと考える。

謝辞

本研究開発の機会を与えて下さり、討論に加わっていただいた当社関係各位に深く感謝の意を表します。

参考文献

- [1]伊藤, 高橋: データベース用自然言語インタフェースにおけるデータベースの自動連結, 第38回全国大会講演論文集

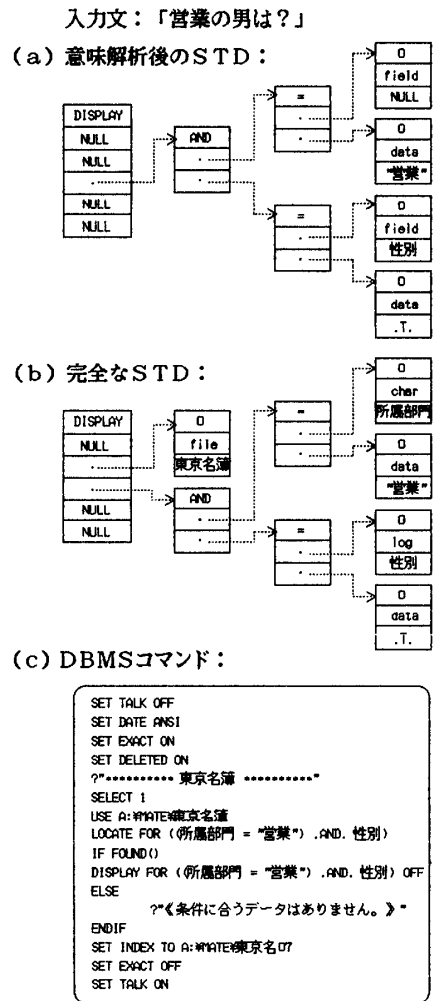


図3 STDの例