

# 文書（テキスト・データ）のデータベース化についての研究

## —— パソコン通信ログ・データベースを例として ——

真田 克彦\*, 山下 陸夫\*

(1993年10月15日 受理)

Investigation of Constructing Database of Documents (Text Data)  
As an Example, Database with Log of Telecomputing

Katsuhiko SANADA, Mutsuo YAMASHITA

### 研究の目的と概要

近年ワープロの普及はめざましいものがあり、それを利用して文書がどんどん作られ、ディスクに電子情報として保存されている。さらに、ネットワークの発達と普及により、それを利用した文書の交換も盛んに行われている。このようにして蓄積される文書の量は、個人においても、組織においても、膨大なものとなってきている。そのため、蓄積された文書をデータベースとして保存しておき、いつでも必要な文書を検索して取り出せるようにしておくことは、今後非常に重要な問題となるであろう。

しかしながら、従来のリレーショナル・データベースは、必ずしも文書の蓄積と管理を得意とはしていない。そのデータ構造はリレーションという一種の表構造に限定されており、文書のような可変長で構造が一定でなく複雑なデータを扱うには適していない。

さらに、従来のデータベースの特徴は、汎用的なデータをできるだけ多くの利用者に共有してもらうことにより、データベースの経済性を追求してきた。データが多くの利用者に共有されるために、高い頻度でのデータの更新処理が要求され、その時点での最新の情報を管理することが第一の使命とされてきた。このようなリレーショナル・データベースの基本思想もまた、文書をデータとして保存することと相容れないものと考えられる。

文書のデータベースは、文書作成や研究・調査など、創造的かつ少人数による作業を支援することが想定される。また、データの更新の頻度はそう高くなく、過去のデータを繰り返し参照することも多い。しかし、更新作業では、大量のデータの処理が必要である。

---

\* 鹿児島大学教育学部数学科

\*\* 川内職業能力開発短期大学校

また、従来のデータベースは汎用コンピュータによる集中処理を想定しているが、文書のデータベースでは、少人数の利用を想定しており、ワークステーションやパーソナルコンピュータによる分散処理が適している。

本論文では、第1部で、リレーショナル・データベース・システムによる文書（テキスト・データ）のデータベースを構築する方法について検討するために桐データベース・システム<sup>\*</sup>を用いて、パソコン通信のログ(log)のデータベースを構築し、その方法や問題点などを明らかにした。

第2部では、最近注目されているオブジェクト指向データベースは、複雑なデータ構造を扱うことを目的としており、分散処理アーキテクチャに適している。このような特徴が、文書のデータベース化に適していると思われるので、オブジェクト指向データベースによる文書のデータベース化の可能性について検討した。

## 第1部 パソコン通信ログファイルのデータベース構築

### 1 はじめに

パソコン通信人口も既に150万人を突破し、双方向のニューメディアとして定着している。これにともなって各通信センターに書き込まれたメッセージの量は膨大なものになっている。この蓄積されている情報を有効に利用する手段や方法の開発が急がれる。

NIFTYのプログラム言語フォーラム(FPL)のBASICの部屋のデータライブラリのログファイルを利用してデータベースの作成を検討中である。現在までにログファイルを利用したデータベースの3つのプロトタイプを作成した。パソコン通信でダウンロードしたファイルから、jgawkを利用してファイル変換処理を行い、約2500個のタイトルリストデータベースを短期間に作成できた。

### 2 データベース作成の基本的な考え方と使用パソコン

#### 2.1 ログファイルによるデータベース作成の意義

パソコン通信で蓄えられたメッセージのデータベース化には次のような意義が考えられる。

- ①パソコン通信のメッセージは増加する一方であり、利用方法の開発による効果が期待できる。
- ②メッセージはコンピュータ入力済みであり、入力なしにコンピュータ利用可能なデータである。
- ③対象にした電子会議室のメッセージは、Q&A的なやり取りの中での情報であり、教育的に利用可能な多くの情報が含まれている。これらの情報を有効活用する方法を開発することは有益である。
- ④情報社会でデータベースを作成することは、個人レベルでもコンピュータの活用を図る観点からも大切であると考えられ、データベース作成実習のカリキュラム開発や社会人を対象にした公開講座に利用可能である。

---

<sup>\*</sup>桐データベース・システムは、管理工学研究所の商標である。

## 2.2 データベース作成の考え方

情報化の進展により個人レベルでもデータベースを作成して、活用することが望まれる。身近で、手軽に入手可能な情報の一つにパソコン通信で得られるログファイルがある。

作成にあたっての基本的な考え方は次の通りである。

- ①当面はデータベース用のソフトは作成することなく、普及している市販のソフトを利用する。  
(将来は、簡単なものを作成することも考える必要がある。)
- ②ハードウェアについても汎用のパソコンを対象とする。
- ③作成手順を単純化して、出来るだけ自動化を目標にする。このためにテキストファイル処理を行う sed, jgawk 等のフリーソフトウェアを活用する。
- ④データベースが作成できたら、キーワード等で該当するメッセージを検索した後、メッセージが読み出せる機能を持たせるように考える。

## 2.3 作成に使用したパソコンシステム

作成の考え方に基づいて、使用したハード及びソフトの概要は次の通りである。

利用対象パソコン: PC-9801シリーズ

(実際には桐を利用できる PC-9801NS/E を使用)

データベースソフト: 国内で開発・市販されている「桐」(Ver 4.0)

その他の利用ソフト: テキストファイル変換等に利用

フリーソフトウェアを主体とした基本ソフト

FD (ファイル管理ソフト)

TED (エディタ)

WXP (仮名漢字変換のフロントエンドプロセッサ)

sed 及び jgawk

使用システムは、身近にある利用可能なハードウェア及びソフトウェアを利用している。データベースソフトの選定にあたっては特別な選定理由は存在しない。

## 3 パソコン通信ログファイルによるデータベースの作成

### 3.1 作成の目的

蓄積された情報の有効利用のためのデータベースの作成方法を確認し、その活用法を検討する。当面の対象は NIFTY の FPL のデータ資料 (BASIC の部屋) によるプロトタイプ of データベースの作成を目標とする。作成対象としたログファイルは BASIC に関する電子会議録であるので、BASIC に関する学習支援データベースとして利用可能であると考えられる。

### 3.2 NIFTY のログファイルとデータベースの読み込みファイル形式

NIFTY の電子会議室のログの形式は図 1-1 のようになっており、データベースソフト「桐」が読み込み可能なファイル種別は「テキスト」と「K3 フォーマット」である。これらの形式のデー

タのサンプル例を図1-2に示している。ログファイルをこれらの形式に変換する必要がある。今回は変換作業が容易なテキスト形式を採用した。図1-3に変換後のテキストファイルの例を示している。

パソコン通信で得られたログをもとに、データベースの項目等に合わせて、読み込みデータのテキストファイルの変換処理を行う必要がある。該当する項目にデータがない場合でもデータを読み込んだとき項目がずれないようにコンマ(,)で区切っておくことが大切である。

```
001/299 SD100379 きだ あきら この会議室は
( 3) 93/03/22 01:48 コメント数:1
BASIC 言語と、そのプログラミングに関連した話題を取り扱います。
Pascal にも負けず、C にも誘惑されず、Modula-2 の厳格さも、perl の
チャランポランさも気にせず、ひたすら最高の道を歩んでいた BASIC も、
最近になって Visual BASIC という新参者とイベントドリブンという新しい
スタイルによる改革の波が迫りつつあるようですが...
とにかく BASIC の話は、この部屋までどうぞ。
SD100379 (FPL) きだ あきら
電子会議 (1:発言 2:コメントを読む 改行のみ:読む) 通常モード
>
002/299 TAB01427 Harry RE:この会議室は
( 3) 93/03/23 21:04 001へのコメント
ってわけで、
何とWindowsの動く環境もなくてVisual basicを買ってしまいました。
さあ、動かしんを買って!
```

図1-1 NIFTYの電子会議室のログ

```
001.tak,データベースの利用法,"001","tak","データベースの利用法",
002.yana,データベースの説明,001,"002","yana","データベースの説明","001"
003.yuki,データベースの入力法,002,"003","yuki","データベースの入力法","002"
004.ame,データ活用法,001,"004","ame","データ活用法","001"
005.mika,データの再利用,001,"005","mika","データの再利用","001"
006.yana,データベースの説明(2),002,"006","yana","データベースの説明(2)","002"
007.yuki,データベースの入力法(2),003,"007","yuki","データベースの入力法(2)","003"
008.ame,データ活用法(2),004,"008","ame","データ活用法(2)","004"
009.yana,データベースの説明(3),006,"009","yana","データベースの説明(3)","006"
010.yuki,データベースの入力法(3),007,"010","yuki","データベースの入力法(3)","007"
011.ame,データ活用法(3),008,"011","ame","データ活用法(3)","008"
```

a) テキスト形式 b) K3形式

図1-2 データベース読み込み可能ファイル形式

```
001/299 SD100379,きだ あきら,この会議室は,( 3),93/03/22 01:48,コメント数:1,BAS
IC,プログラミング,Visual BASIC
002/299,TAB01427,Harry,RE:この会議室は,( 3),93/03/23 21:04,001へのコメント,Visua
l basic,
003/299,PPA03560,録音機,はじめましてのあいさつ,( 3),93/03/23 22:13,コメント数
:1,True BASIC,Back to BASIC,ホタテマス大のゲームとカー
ン
```

図1-3 変換後の図1のテキストファイル

### 3.3 データベース作成手順と項目の決定

パソコン通信のログファイルを利用したデータベースの作成手順は次の通りである。データ入力の手数は省略できるが、データの変換作業が問題となる。

- ア) データベースの項目の決定
- イ) データベース用ファイルへの加工・編集
- ウ) データの読込とデータベースの作成
- エ) 作成データベースの試用と評価

#### (1) データベース項目の決定

図1-1に示すログを参考にすると、データベースの項目の候補としては次の項目が挙げられる。

- ア) メッセージ番号, イ) ID, ウ) ハンドルネーム, エ) タイトル, オ) 登録日時,
  - カ) 会議室番号, キ) 関連メッセージ番号, ク) コメント数, ケ) キーワード, コ) メッセージ
- ア)~コ)の項目から必要に応じて適切な項目を選定してプロトタイプのデータベースの項目として定めた。試作した各プロトタイプの項目は次の通りである。

- プロトタイプ1 上の項目の候補からメッセージを除いた項目で作成
- プロトタイプ2 上の項目の候補からキーワードとメッセージを除いた項目で作成
- プロトタイプ3 上の項目の候補からキーワードを除いた項目で作成

(実際の作成はタイプ2にメッセージを結合して作成している)

キーワードは、メッセージを実際に読んで本文中から適当な語句を選定してキーワードとしており、その選定にはかなりの時間と労力を必要とする。今後は、メッセージとキーワードとの関係を更に検討する必要がある。また、メッセージの長さは一定ではなく、長いメッセージは分割して

データベース化する必要がある。データベース項目の決定にはメッセージの長さを考慮する必要がある。

## (2) メッセージ番号の付け方

NIFTY (FPL) の BASIC の部屋のこれまでのメッセージは約500個毎に整理され、データライブラリに格納されている。現在、既に約2500個のメッセージが5回に分けられて分割保存されている。保存されたメッセージは各回とも001から約500までの番号が付けられており、データベース化にあたっては、次の例に示すように格納回数を上位1桁に1～5の番号で付加して全体のメッセージ番号を付けている。メッセージを項目に持つプロトタイプ3では、メッセージが長いものはメッセージを分割して作成したので、その分割番号を下位に付加している。

例 3302 3回目に格納された302番のメッセージ

11502 1回目に格納された150番のメッセージで分割された2番目のメッセージ

## 3.4 読み込みファイルの加工・編集

使用したデータベースソフト「桐」では、1レコードは1000文字以内と制限されている。メッセージを項目に入れる場合にはこの点を考慮すべきである。以下に各プロトタイプのデータベース用読込ファイルの編集・加工について述べる。

### (1) プロトタイプ1の場合

ログファイルをエディタ (TED) に読み込み、手作業で区切り記号の挿入及びキーワードの選定を行い、読込用ファイルを編集・加工した。編集作業及びキーワードの選定を手作業で行うにはかなりの労力が必要である。

### (2) プロトタイプ2の場合

sed でタイトルリスト一覧の部分を切り出し、区切り記号の半角のスペースをコンマ (,) に変換し、更にメッセージ番号に分割のブロック番号を付け加える作業を jgawk を利用して行って読み込みファイルの編集・加工を行った。項目で空白の場合には項目がずれないように区切り記号を必ず挿入するように注意が必要である。

### (3) プロトタイプ3の場合

ログファイルからメッセージ番号とメッセージのみを切り出して、メッセージ本文を1項目となるように擬似改行コードを挿入して編集作業を行う。メッセージが長くて1項目の範囲を超える場合は分割してメッセージ番号の最下位に分割番号を付加する。メッセージ番号とメッセージのデータベース表を作成し、プロトタイプ2と結合してプロトタイプ3を作成した。

メッセージの編集は現段階では手作業で行っており、エディタの編集機能を利用して「擬似改行コード」を挿入している。メッセージの編集を手作業で行うにはかなりの労力が必要なので、今後編集作業を改善する余地がある。

ここで「擬似改行」とは、エディタなどで1行の長さに制限がある場合に、擬似的な改行コードを加えることによって1行とみなすようにすることを言う。使用した TED の場合には“CTRL”



ファイルから簡単にデータベースの読み込み可能ファイルに変換できる方法を開発すべきである。

#### 4 試作データベースの概要

現在, 3つのプロトタイプ of データベースを試作して評価を行った結果, 最終的にはプロトタイプ3にキーワードを付加したものを作成する方が望ましいと考えている。以下に作成した3つのタイプの概要を述べる。

##### (1) プロトタイプ1

キーワードの項目を持っていることが特徴である。しかし, 初めて作成したタイプであることもあり, 図1-8に示すようにデータに冗長がある。また, キーワードの選定にかなりの労力を要するので, メッセージの項目を持つプロトタイプ3で, データベースの機能を利用して, 直接メッセージを検索する方法が, 容易にデータベースの作成ができるものと考えられる。ハンドルネームや会議室番号は省略可能である。

:メッセージ:	ID	:ハンドル:	タイトル	:会議室:	日時	:開催:	コメント数	:	キーワード
001/239	SD100379	紀 敏	この会議室は	( 3)	93/03/22 0		コメント数: 1		BASIC プログラミング Visu -BASIC
002/239	TAB01427	Harry	RE:この会議室は	( 3)	93/03/23 2 001^				Visual-basic
003/239	PFA03560	藤 提督	はじめましてのごあいさつ	( 3)	93/03/23 2		コメント数: 1		TrueBASIC Back-to-BASIC * トラスのクエリ
004/239	TAB01427	Harry	QBについての質問	( 3)	93/03/24 1		コメント数: 1		QBEOF QBユーザー定数型の引渡し
005/239	JAD00247	J.T.	QBのEOFについて	( 3)	93/03/25 1 004^		コメント数: 2		QBEOF QBでは連続レコードでEOFは 次のレコードGETした
006/239	HGD03067	じゃん	RE:QBのEOFについて	( 3)	93/03/25 2 005^				シーケンスファイルのEOFは最後のデータを数バイトで
007/239	TAB01427	Harry	RE:QBのEOFについて	( 3)	93/03/26 0 005^		コメント数: 1		QB C EOF
008/239	SD100379	紀 敏	RE:QBのEOFについて	( 3)	93/03/26 2 007^		コメント数: 1		BASIC QB EOF シーケンスファイル ランダムアクセス EO
009/239	HGD00201	英斗恋	EL_BASIC ver.1.53Tの不具合について	( 3)	93/03/27 0				EL_BASIC 登録
010/239	TAB01427	Harry	RE:QBのEOFについて	( 3)	93/03/27 0 008^				QB 書き方
011/239	NBA00671	Nydel	これはこれは	( 3)	93/03/27 1				新しい修飾の機能 VBDOS イベント・トリップ
012/239	QCA01641	きゅろ	VBがきました, けど(…)	( 3)	93/03/27 1		コメント数: 1		VB -for-DOS VB のハング
013/239	TAB01427	Harry	RE:VBがきました, けど(…)	( 3)	93/03/27 1 012^		コメント数: 2		VB
014/239	QCA01641	きゅろ	VXはだめ?!	( 3)	93/03/28 1 013^				DA VX 286 VM VBDOS
015/239	QCA01641	きゅろ	原因はわかった	( 3)	93/03/29 1 013^				VBDOS メモリドライバの性能 NEC検証のEMM386.SYS
016/239	HBA00270	Still	初めまして! さっそく質問なんですが	( 3)	93/03/29 1		コメント数: 1		QBのPalette文 QuickBasic
017/239	HGD02533	んたんび	86BC、注文しました	( 3)	93/03/29 2				FA(factory automation)のシステム設計 QB VB
018/239	HBA00270	Still	RE:初めまして! さっそく質問なんですが	( 3)	93/03/29 2 016^		コメント数: 2		Cで使うARGV(0)の(自分自身の実行ファイル名をフルパスで取得)QB
019/239	GFC01740	あちゃ	BASICの定義って…?	( 3)	93/03/31 0		コメント数: 2		QuickBASIC TrueBASIC N88-BASIC X68000
020/239	QFF03453	坂田 肇	QBの質問/文字列の文字抽出方法?	( 3)	93/03/31 0		コメント数: 2		QBの決まった文字列の抽出 文字列が入れられた変数ASから、半角英数字の
021/239			メッセージなし						
022/239	TAB01427	Harry	RE:QBの質問/文字列の文字抽出方法?	( 3)	93/03/31 2 020^				QB CSN\$関数
023/239	NAC03056	玉井 孝	届いたぞ	( 3)	93/04/01 2				T K W - 8 6 B C V e r 5 . 5 2
024/239	SD100379	紀 敏	RE:初めまして! さっそく質問なんですが	( 3)	93/04/02 0 018^		コメント数: 1		COargv[0]に相当するものを得るプログラム
025/239	SD100379	紀 敏	RE:BASICの定義って…?	( 3)	93/04/02 0 019^		コメント数: 2		BASIC の定義 名BASIC共通している点, 代入記号'=, 修飾的 変数した他言語に似て, 分かりやすいように作られたと入る意味でBA
026/239	GFC01740	あちゃ	RE:BASICの定義って…?	( 3)	93/04/02 1 025^				
027/239	HBA00270	Still	発言者削除	( 3)	93/04/02 1 025^				
028/239	HBA00270	Still	RE:初めまして! さっそく質問なんですが	( 3)	93/04/02 1 024^				統合開発QB.EXEを実行プログラムではその名前を尋ねる事ができました
029/239	PDD01714	とりちゃん	86BC ver 5.52 ...	( 3)	93/04/02 2				
030/239	PDD01714	とりちゃん	BASICとは?	( 3)	93/04/02 2		コメント数: 2		文字変数を***\$で扱う PRINT FOR IF GOTO GOSUB
031/239	HGD03067	じゃん	RE:BASICの定義って…?	( 3)	93/04/03 0 019^				ベータは, 数値用にはきたも 変数を宣言しなくても使えるのはベータ
032/239	NBA00671	Nydel	コマンド文字列の取得	( 3)	93/04/03 0 018^		コメント数: 2		MS-DOSの3.3かそれ以降であれば, 入力した文字列は環境変数環境に保
033/239	HBA00270	Still	RE:コマンド文字列の取得	( 3)	93/04/03 1 032^		コメント数: 1		PSP0セグメントアドレス取得 int21-62HはDOS-Ver-3.0
034/239	PAG03057	KITA	RE:コマンド文字列の取得	( 3)	93/04/04 0 033^				コマンドライン引数 QB45の標準関数 >QB/LGENで取得
035/239	PDD01714	とりちゃん	re:int21 62H 連続的に実行プログラム	( 3)	93/04/04 0				comファイル処理時のコマンドセグメント管理PSP.exeファイル処理時の
036/239	SD100379	紀 敏	発言者削除 :誤字修正	( 3)	93/04/04 0 030^		コメント数: 1		
037/239	SD100379	紀 敏	RE:BASICとは?	( 3)	93/04/04 0 030^				文字変数を***\$で扱うことはいないBASIC PRINT FOR IF
038/239	SD100379	紀 敏	RE:コマンド文字列の取得	( 3)	93/04/04 0 032^				DOSの3.XXと2.XX #24と同じ処理cargv[0]を*tl
039/239	QCB01521	ockeghem	再びBASICの定義	( 3)	93/04/04 1 036^				Visual-Basic CにはINPUTがありません Ver2.0以降の
040/239	FFA00237	Fling Rock	RE:BASICの定義?について	( 3)	93/04/05 0				BASASICO明解 Beginner's All-purpose Sym
041/239	HGD00201	英斗恋	EL_BASICテスト公開期間の終了について	( 3)	93/04/05 2		コメント数: 1		
042/239	HGD00201	英斗恋	RE:EL_BASICテスト公開期間の終了について	( 3)	93/04/05 2 041^				ver.1.54 複数のCONST変数を一様に宣言できない 文法エラー
043/239	NBF01226	L a y	ようこそ 藤提督さん	( 3)	93/04/06 0 003^				True Basic
044/239	GAA01062	快人	BASICでオプション取得・・・ 快人	( 3)	93/04/06 1		コメント数: 1		MS-BASIC QBでのオプションの取得する方法
045/239	GBF00535	Doc Hollid	RE:BASICでオプション取得・・・ 快人	( 3)	93/04/06 1 044^		コメント数: 1		COMMANDSで取得できます
046/239	GBF00535	Doc Hollid	RE:BASICでオプション取得・・・ 快人	( 3)	93/04/06 1 045^				QBでの話です MSは知りません
047/239	HBA00270	Still	文字列の切り出しについて	( 3)	93/04/10 0		コメント数: 1		数字を含んだ文字列変数 バイト単位で格納時に切り分ける方法 PSP セグ
048/239	JAH02230	げんちゃん	QBとFEPについて, 教えてください。	( 3)	93/04/11 0		コメント数: 1		QB自作のソフト 入力時に自動的にFEPの起動 FEPHATON?!
049/239	TAB01427	Harry	RE:QBとFEPについて, 教えてください。	( 3)	93/04/11 1 048^		コメント数: 1		拡張ライブラリ ファンクション
050/239	TAB01427	Harry	QBでいきなり暴走	( 3)	93/04/12 1				QB 巨大プログラム

図1-8 プロトタイプ1の表形式

##### (2) プロトタイプ2

ログファイルの中のタイトルリストのファイルを利用して, jgawk を使用し, 約2500件のメッセージのタイトルデータベースがある。ID やタイトルで検索等が可能で, 作成したプロトタイプ

