社会調査結果の視覚化データベース

Visualized Database of Social Survey Results

吉田 光雄

Mitsuo YOSHIDA 大阪大学 人間科学部 Human Sciences, Osaka University

キーワード: Mathematica, 探索的データ解析, 情 報処理教育

Keywords: Mathematica, EDA(Exploratory Data Analysis), Education for Information Processing

あらまし:社会調査結果を計算機内に保存し、必要に応じてデータを取り出し、Mathematicaを用いて、基本統計量を計算すると同時に、ヒストグラム、散布図等のグラフを描き、データの様相を視覚的に検討することのできるデータベースをワークステーション上に構築した。すなわち、保存されたデータをテキストとして取り出せると同時に、それらを視覚的に提示し、探索的にデータの内部構造を探り、データの持つ情報を十分に引き出して分析を深化させるために用いることができる。

きる。 本システムは小規模の試作用データベースで あるが、メモリー・ハードディスクの容量や高速 演算などの十分な計算機資源が得られれば、さら に大規模のデータにも適用することができる。快 適に動かして大量のデータを処理するためには、 かなりな量のメモリを必要とするが、統計処理 のメニュの提示や選択、さらにはより美しいグラ フィックスの描画や3次元表示法に対する工夫な ど、ユーザー・インターフェイスの改良は今後の 課題である。

Summary: Social survey results should be preserved as database, so that one can consult it afterwards and extend it to another analysis. If it can be browsed in a visual format, such a database would be an efficient and convenient tool for many kinds of social surveys.

This report summarized how to save data in Mathematica as a database, and how to make statistical treatments using graphic techniques, such as histogram, bar-chart, pie-chart, scatter diagram, three dimensional graphic and others in it.

Mathematica seems to be suitable for this purpose, because of its wide coverage of mathematical and statistical usages.

This database seems to be useful, not only to preserve all of the data, but as a visual database that contributes as an EDA(Exploratory Data Analysis) and VDA(Visual Data Analysis) revealing new and hidden structures of the data.

An example of social survey results concerning computer education at Osaka University was demonstrated.

1. はじめに

Mathematica ^{TM(3)}は 1986 年、Wolfram Research 社から発売され、現在世界中で広く使用さ れている数式処理ソフトのひとつである。操作の 容易さ、機種互換、グラフィックス化等に特徴が あり、例えば、式の展開、因数分解、方程式、数 値計算、行列演算、微分・積分、微分方程式など の数式処理の他、関数のグラフ表現、統計計算な ど、数学的処理を容易に行うことができる。強力 なパッケージ⁽²⁾を多数含み、"A System for Doing Mathematics by Computer"として評判の高 いものである。

統計学に関するパッケージも多く含まれている ので、さまざまな統計処理にも活用することがで きる。しかも、グラフィックス機能が豊富である ので、目で見る統計学として有効である。ワーク ステーション (NeXT) 版およびマッキントシュ版 はユーザー・インターフェイスのよい Notebook を介して使用可能であり、コマンド入力、プログ ラミング、出力、グラフィックスがすべて同一ウ インドウ上で実行され、しかも両機種間で高い互 換性を有している。

統計パッケージ SAS を用いた社会調査結果の 視覚化データベース、ならびにS言語を用いたも のも可能であるが⁽⁵⁾、本稿では Mathematica を 用いて行う方法について報告する¹。また、まだ モジュール化されていない、いくつかの統計的手 法については、内蔵の言語を用いてプログラミン グを行ったので、それらについても報告する。追 加された各種の統計演算はライブラリとして保存 されている。

Turkey,J.W.⁽⁴⁾により探索的データ解析 (EDA, Exploratory Data Analysis) の方法が言われて以 来、統計データを視覚的に分析する方法 (VDA,

¹本稿は、文部省科学研究費補助金・重点領域研究、「情報 化社会と人間」研究成果報告書、および日本計算機統計学会 第8会大会にて報告したものの一部である。

Visual Data Analysis) が注目され、いろいろな に合わせて並べ換えておく。 方法が提案されている。本データベースはあくま でもデータの統一的保存が目的であるが、保存さ データを一括して読み込ませておくと迅速な処理 れたデータをテキストとして取り出せると同時に が可能であるが、十分でない時は、必要な項目を 視覚的に提示し、全体像の把握を容易にすること 処理の都度読み込めばよい。 をも目指す。そして、更に必要な場合には再度統 計処理も行うことが可能なものとする。統計処理 はすでになされているが、様々な角度からより探 索的にデータの内部構造を探り、データの持つ情 報を十分に引き出し、分析を深化させるために用 いることができる。

大阪大学では、情報処理教育センターを中心に NeXT を端末とする分散処理システムのネット ワークが構築され、学生の情報処理教育に活用さ れている(1)。教材の提供・課題の提示・レポート の収集などが容易に行える教育支援のためのアプ リケーションが開発されており、それを活用する ことにより、登録学生は容易に本データベースに アクセスすることができる。

2.1 保存および参照

Mathematica を起動する前に、カレント・ディ レクトリにデータをテキスト (Ascii) 形式で保存 しておく。データの形式は、各行にサンプル(ケー スまたはオブザーベーション)、各列に変数(調査 ためには、統計パッケージを予め読み込んでおか 項目)をあて、行列データの入力形式に従って、 ねばならない⁽²⁾が、これらの一連の初期作業を 各行ごとに { } を使用する。

データでもよい。ただし、処理や結果の解釈に際 に保存されているので、それらを利用して copy して、データの属性に注意する必要があり、現在 & paste しつつ、再度実行することも可能である。 のところデータベースのシステム内に、可能な統 計処理をチェックする機能は組み込まれていない。

調査票の項目一覧も必要であり、Mathematica 内の変数番号と、調査票の項目番号の対照も必要 である。本計算機システムが NeXT ワークステー 可能であり、項目対照を明示するウインドウを常 時開いておくと、操作が容易となる。項目の選択 は変数名やラベルを使用するのではなく、通し番 号で行う。

調査票は多くの場合、回答者のバイアスを避け るためランダムに並べられているが、計算機に保 や母平均に関する区間推定・検定を行うことができ 存する際には、予め質問項目番号、回答カテゴリ の方向性などを整理し、整合性を保って保存して をベクトルとして取り出す。その際、オプション おいた方が便利であろう。調査項目一覧もデータを用いて、有意水準、両側・片側などを設定する。

メモリなど、計算機資源が豊富な場合には、全

一旦保存されたデータは修正ができないよう、 ライトプロテクトをかけておく。データを加工し ても、保存はせず、必要な修正はその都度行うよ うにしておかないと、データベースの機能を果た さなくなる。従って、以後の修正が無用となるよ う、最初に計画的に保存しておく必要があろう。

現在のところ NeXT 版 Mathematica (Ver. 2.1) は日本語対応とはなっておらず、出力に日本語が 使用できないのは難点である。ただし、処理結果 をワープロや作図ソフトに取り出して編集すれば、 日本語は使用可能である。

2.2 Mathematica の基本

Mathematica を起動し、ハードディスク上の カレント・ディレクトリ (~) に保存されている **2**. データの作成と Mathematica の基本 データ・ファイルを Mathematica に読み込む。 必要に応じて、特定の項目に関する全サンプルの 回答を抽出し、編集しなければならないが、そう した作業を便利に行うことのできるコマンドも多 く用意されている(表1)。また、統計処理を行う ルーチン化してプログラミングしておいてもよい データは連続量でもよいし、数字化された属性 し、 Mathematica との対話の内容は Notebook

3. 統計処理

3.1 基本統計量の算出

統計パッケージ(記述統計学)を用いることに より、平均・メディアン・モード・分位数などの ションであるため、多重ウインドウを開くことが 代表値、分散・不偏分散・標準偏差・範囲・四分 偏差などの散布度のほか、歪度・尖度なども直ち に計算することができる(表3)。

3.2 区間推定·検定

パッケージ(信頼区間、検定)を用いて、母分散 る(表4)。データは項目番号を指定して粗データ

表 1: データ読み込み・編集のため	のコマンド
--------------------	-------

	パッケージ Statistics'Master'		
	Statistics `DataManipulation`		
	ReadList["~/file1.dat",Number]		
	ReadList[" \sim /file1.dat",Number,		
RecordList->True]			
	TableForm[data]	Column[data, n]	
	$Column[d, {n1, \cdots}]$	ColumnTake[d,n]	
	$ColumnDrop[d, {n1, \cdots}]$	ColumnJoin[d1,d2,]	
	Row[data, n]	RowTake[d,n]	
	Row[d, {n1,}]	RowDrop[d, $\{n1, \cdots\}$]	
	RowJoin[d1,d2,]		
	BooleanSelect[d,sel]	TakeWhile[d,pred]	
	LengthWhile[d,pred]		
1			

表 4: 推定・検定のためのコマンド
パッケージ Statistics'Confidenceintervals'
MeanCI[data]
MeanCI[d, KnownVariance->v]
MeanDifferenceCI[d1,d2]
MeanDifferenceCI[d1,d2,
<pre>KnownVariance -> {v1,v2}]</pre>
VarianceCI[d]
VarianceRatioCI[d1,d2]
パッケージ Statistics'HypothesisTests'
MeanTest[d,mu]
MeanTest[d,mu, KnownVariance->var]
MeanDifferenceTest[d1,d2,diff]
VarianceTest[d1,var]
VarianceRatioTest[d1,d2]

表 2: データ集計のためのコマンド

パッケージ Statistic	s'DataManipulation'	
Frequencies[data]	QuantileForm[d]	
CumulativeSums[d]	BinCount[d, {min, max, dx}]	
RangeCount[d, {c1,c2,dx}]		
CategoryCounts[d, {cat1, cat2,}]		
BinLists[d, {min, max, dx}]		
RangeLists[d, {c1, c2, dx}]		
CategoryLists[d,{ca	t1,cat2,}]	

表 3: 基礎統計量算出のためのコマンド

Statistics' DescriptiveStatistics'

Median[d]

GeometricMean[d]

パッケージ

Mean[data]

Mode[d]

パッケージ (連続型分布) には、予め正規分布、 *t*-分布、*F*-分布、*χ*²-分布 などの標本分布が組 み込まれているので、上側・下側・両側確率、上 側・下側・両側パーセント点を求めることができ る(表 5)。従って、任意の値の自由度についての 三統計数値が計算可能であるし、プログラミングす ることにより、粗データからではなく、標本平均・ 標本分散などから、直接区間推定や検定を行うこ ともできる。

4. データの視覚化

グラフ用パッケージを読み込み、グラフを描く ことができる (表 6)。BarChart で描くのは棒グ ラフであるが、ヒストグラムの代用とすることも できる。

標本統計量から分布のパラメータを求め、確率 分布をあてはめて、密度関数・分布関数などのグ ラフを描くこともできるし、データの分布と重ね

Quantile[d,q]	
LocationReport[d]	表 5: 標本分布
VarianceMLE[d]	パッケージ Statistics'ContinuousDistributions'
MeanDeviation[d]	NormalDistribution[0, 1]
SampleRange[d]	ChiSquareDistribution[df]
Skewness[d]	StudentTDistribution[df]
Kurtosis[d]	FRatioDistribution[df1,df2]
ShapeReport[d]	NonCentralChiSquareDistribution[df,lambda]
correlation[x1, x2]	NonCentralStudentTDistribution[df,lambda]
partialCorr[x, i, j]	NonCentralFRatioDistribution[df1,df2,lambda]
contTable[x1, x2]	CDF[dist,x]
	Quantile[dist,q]
	Quantile[d,q] LocationReport[d] VarianceMLE[d] MeanDeviation[d] SampleRange[d] Skewness[d] Kurtosis[d] ShapeReport[d] correlation[x1,x2] partialCorr[x,i,j] contTable[x1,x2]

表 6: グラッフィクス

cs'		
$BarChart[d1, d2, \cdots]$		
t [d]		
ListPlot[d]		
DisplayTogether[plot1,plot2,]		
DisplayTogetherArray[{{p1,p2,},		
$, \dots, \{\dots\}\}$		
cs3D'		
Plot3D[d]		
eListPlot'		
MultipleListPlot[d1,d2,]		
$MultipleListPlot[d1, d2, \cdots, PlotJoined->True]$		

て描き(表 6)、確率分布のあてはめの程度を検討 することもできる。

調査項目 (変数) をいろいろ抜き出し、グラフ を描くことによって、データ全体を見渡すことが 出来るし、かつ outlier などの発見も容易となる。

グラフィックス・コマンドはオプションが豊富 に用意されており、凡例・カラー・線の太さと種 類・ラベル・テキスト挿入などを自由に設定するこ とにより、様々な美しいグラフを作成し、プレゼ ンテーションの効果を上げると同時に、EDA と しても活用することができる。データから探索的 にさらに多くの情報を引き出すことが、本データ ベースの目的でもある。

5. ユーザー・グラフィックス関数

データの種類、例えば、連続量データかカテゴ リデータかにより、描くグラフの種類も異なって くる。どのような種類のデータに対して、どのよ うなグラフが描けるかを、先の標準パッケージを もとに、描画し易いようにユーザー関数としてプ ログラミングを行った。

社会調査の場合、属性データの分類に際しては、 0,1,2,・・・の離散型データが用いられ、名義尺度と して数量化されるケースが多い。さらに順序デー タ、連続量データであっても離散型で処理される ケースも多く、属性データの処理は基本であろう。

連続量データの場合でも度数分布表に集計して、 順序尺度として使用されるケースが多い。従って、 多用されるのも、カテゴリカルデータのグラフ化 である。

統計処理としては、度数の集計、度数分布、分

表 7: ユーザー・グ	ラフィックス関数
1 変数データの統計	十図
11.msd[dname, item]	平均・標準偏差
12 myec[vec]	平均値ベクトル

12.mvec[vec]	平均値ベクトル	
13.psing[vec]	比率	
14.distsing[vec]	度数分布	
2 変数データの統計図		
21scatter[vec]	散布図	
22.ctab[vec]	相関表	
23.cortab[vec]	相関表 (多数)	
多数データの統計図比較		
31.hist6[vec]	ヒストグラム	
32.mvprof6[vec]	平均値プロフィール	
33.pcomp6[vec]	比率の比較	
34.distcomp6[vec]	度数分布の比較	
35.mpc6[vec]	平均 (比率) の比較	

割 (関連) 表などであり、そららの図示として、円 グラフ、帯グラフ、分割表の 3D 表示などが用い られる。

連続量データの全体像を見るためには、粗デー タをそのまま用いて、散布図を描くこともできる が、プレゼンテーションを工夫して、3次元(立 体)散布図を描くことも可能である。

さまざまなグラフ化に対応するため、標準グラ フを活用したユーザ関数を作成した。パッケージ を使用しているので、それらの読み込みが必要で ある。

データセットから項目 (変数)を指定して、平 均・標準偏差を算出する関数が msd である。デー タ名を dname で与え、項目を item で選出する。 item が複数に及ぶときは、mvec を使用する。こ のとき関数の引数は vec のみであるので、関数 を呼ぶ前にデータセット dname[1]=datasetname と項目番号 vec=itemnumbers を予め設定してお かねばならない。vec はリストである。mvec は グラフを描く。

データが (0,1) の場合は、回答 (yes=1) の比率 を求め、ヒストグラム (棒グラフ)を描く。

下位反応カテゴリの度数を集計し、比率を求 め、分布状況についてのヒストグラムを描くのが psing である。これらは複数項目の処理が可能で あり、事前にデータセット名と処理項目番号を選 定しておく。

2 変数の連続量データについて、平均・標準偏 差・共分散・相関係数を算出し、散布図を描くの が scatter であり、項目を vec で選択する。選択 処理する。

離散量データの場合、セルの度数のカウントと なり、散布図では度数が1点としてしか図示され ないので不合理である。図示方法を変更し、3次 元立体図と3次元表面図で図示するのが ctab で ある。1 項目の単純集計 (marginal distribution)、 2 変数のクロス集計 (joint distribution) の度数も 出力する。単純集計のグラフでは標準関数の制約 から、2項目のみの描画ができず、従って1項目 を反復して、並べて描くこととする。

ctab は2項目間の処理のみ、cortab は多数項 日を vec で与え、その中から順次2項目の組み合 わせを取り出して処理するものである。

以上は単独のデータセットについての処理であ るが、それが複数となったとき、データ間の比較を 行うためのユーザー関数が、番号 31 以下である。

まず、hist6 は度数分布を集計しヒストグラムを 描く。mvprof6 はで指示する項目の平均値を求め、 それのプロフィールを描く。pcomp6 は比率の比 較、distcomp6 は度数分布の比較である。mpc6 は平均値(または比率)の比較を行う。

予めデータセット名 (dname[i]=filenamei) と その数 (ng≤ 6) を与え、処理する項目を (vec= *i*₁, *i*₂, …) で与える。データセットの数の上限は 6 である。あまり多くのグラフを重ねると相互の 比較が困難となるが、上限の6をプログラム内で 修正することは容易である。

グラフはいずれも折れ線グラフを描き、相互比 較の便を期す。

6. 多変量解析

現在のところ、多くの多変量解析の手法は Mathematica に組み込まれておらず、他の統計ソフト を使用しないのであれば、自らプログラミングせ ざるをえない。行列やベクトルの積、逆行列など の行列演算や、固有値・固有ベクトルの計算など は Mathematica 内に定義されているので、それ らを用いれば多変量解析のプログラムを書くこと は、さほど困難な作業ではない。むしろ、 Mathematica に内蔵のプログラム言語では、ベクトル 演算がサポートされており、一般の言語より短く、 簡単に書くことができる。

これらを発展させ、多変量解析の例として、重 回帰分析、判別分析、主成分分析のプログラミング を行った。必要に応じてデータベースから簡単に

された vec の項目から 2 項目づつを取りだして アクセスすることができる。さらに、FORTRAN、 C. Pascal などの一般の言語で書かれた実行形式の プログラムをリンクさせて、用いることもできる。

7. 具体例

社会調査の一例として、大阪大学生に対して実 施された情報処理教育に関する調査結果(1)を使用 する。調査項目は以下の通りである。

学部・専攻などのフェースシート項目

- 大学入学以前のコンピュータ経験
- 入学後の計算機使用の実態
- 計算機に対する意識
- 情報処理教育に対する意見

全学部1年生のランダム・サンプルとして 999 名 (在籍学生のほぼ 35%(文科系 489 名、理科系 510 名))のデータを入手し、計算機内に保存し、 データベースとした。

調査項目は、ID 番号をも含めて No.1 ~ No.106 である。5段階表示の順序データ、(0-1)表示のカ テゴリデータなどが混在している。

サンプルは、例えば(1)学部別、(2) 文科系·理 科系別、(3) 全調査対象を合併したもの、などの 構成とすることが可能であり、3 群毎にファイル を作成しておけば、必要に応じて目的のデータを 読み込むことができる。

基礎統計量の計算、単純・クロス集計、グラフ 化が可能であるが、一例として、図1~図6に各 種のグラフを示す。

煵 文

[1] 松浦敏雄他 (1993) 教育用計算機システムの運 用と学生の意識、行動計量学、40、17-31.

[2] Phillip Boyland (1992) Technical Report, Guide to Standard Mathematica Packages, Wolfram Research.

[3] Stephan Wolfram (1988) Mathematica, Addison-Wesley, 訳書 (1992) アジソン・ウエスレイ.

[4] Tukey, J.W. (1977) Exploratory Data Analysis, Addison-Wesley.

[5] 吉田光雄 (1993,95) 映像情報データベース の開発, 文部省科学研究費補助金・重点領域研究 「情報化社会と人間」,研究成果報告書.

565 吹田市山田丘 1-2, 大阪大学人間科学部 Tel: 06-879-8051, Fax: 06-879-8054 E-mail: voshida@hus.osaka-u.ac.jp



図 1: Q8-2. コンピュータが好きだ (文科 系、賛成の程度)



図 3: Q8 2(Item1)-20(Item2) のクロス集 計 (文科系)



図 5: Q3. 自宅におけるワープロ(W)・パソ コン(PC)の有無と使用(文(淡)・理(濃))



図 2: Q8-20. 卒業後はコンピュータを使用 する仕事をしたい (文 (実線)・理 (点線)の 比較)



図 4: Q8 2-20 のクロス集計 (理科系)



図 6: 主成分得点散布図 (文科系·F2-F3)