

□ CHERNOFFの顔型グラフ (FACE GRAPH)

人間は人の顔を見る事で容易に個人を識別出来、日常的にもこの顔パターンの識別能力が高く、その微妙な変化を読み取る事にも優れているという観点から出発したものが“顔チャート”である。

CHERNOFFの顔チャートでは、顔を形成する様々な要因それぞれに一つの数値データをあてはめる。この数値データの変化により、様々な顔が形成されてゆく。最近では経済指標の表示等によく用いられている。

このチャートは人間の顔という特別な対象を表示図形として用いている為に、作図時にはレーダーチャート等、他の表示手法にはない注意が必要となる。つまり、人間は顔を見て単に識別するだけでなく、「顔の変化」とある一定の「事象」を相関付けて考えるのが自然だからである。顔が笑っている時、事象としては理想的な事が起こっていると考え、泣いている時や怒っている時は希望にそぐわない事が起こっていると考える。

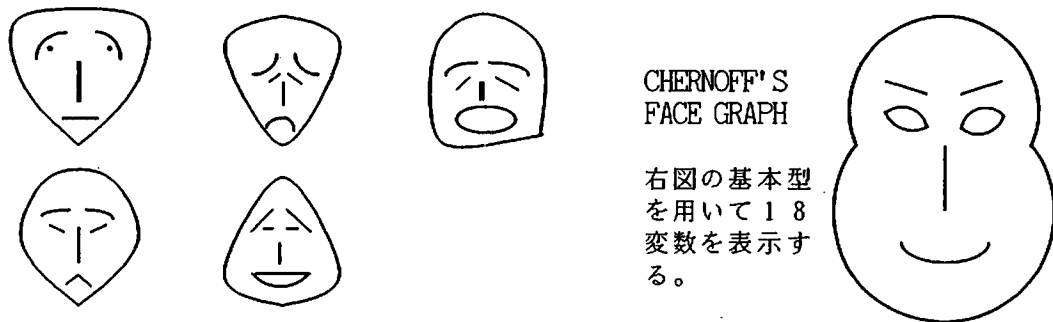


図5-2 顔チャート事例及びCHERNOFFの顔チャート

従ってこのチャートを用いる時は、当面の解析対象が有する意味と顔が持つ意味とが一致するようにパラメータを選択する事が必要である。この選択に失敗すると、例えば経済指標の表示に顔チャートを用いる時、経済が絶好調である時に泣き顔や怒った顔となり、沈滞時に満面の笑みをたたえる事となる。つまり、識別は出来ても解析対象の持つ情報と顔の表情との相関が取れなくなる。このように顔チャートを利用する場合、パラメータの選択は最大の留意事項となる。

この顔チャートと同じ基本概念で構築されたチャートとしては体形グラフ (BODY GRAPH) がある。これは、顔のみならず人間の体形総てを利用してチャートにするものである。

□ CHERNOFFの顔チャートに用いられる顔の作成要因について

顔チャートは顔という限定されたものを対象としている為、使用される数値データの数にも限界がある。CHERNOFFが最初に用いた顔チャート用の数値データ数は以下に示す18種類である。

変数	計算式	変数の意味
X1	$h^* = \{(1 + X1) H\} / 2$	鼻の中心と輪郭迄
X2	$\theta^* = (2 X2 - 1) \pi / 4$	
X3	$h = \{(1 + X3) H\} / 2$	顔の上半分の長さ
X4	X4	上半顔の離心率
X5	X5	下半顔の離心率
X6	X6	鼻の長さ
X7	$p_m = h \{X7 + (1 - X7)X6\}$	口の位置
X8	X8	口の曲率
X9	$a_m = X9(h / X8) \text{ or } X9 W_m$	口の幅
X10	$Y_e = h \{X10 + (1 - X10) X6\}$	目の位置
X11	$X_e = W_e (1 + 2 X_{11}) / 4$	
X12	$\theta = (2 X_{12} - 1) \pi / 5$	目の傾き
X13	X13	目の楕円の離心率
X14	$L_e = X_{14} \min (X_e, W_e - X_e)$	目の幅の半分
X15	X15	瞳の位置
X16	$Y_b = 2 (X_{16} + 0.3) L_e X_{13}$	
X17	$\theta^{**} = \theta + 2 (1 - X_{17}) \pi / 5$	眉の傾き
X18	$L_b = r_e (2 X_{18} + 1) / 2$	眉の長さ

□ 三角多項式グラフ

アンドリュース (Andrews, D. F) により提案されたものである。個々のパターンは $-\pi \sim +\pi$ 迄の範囲で動く一本の曲線として表現される。曲線の式は以下の式を用いて算出される。

・変数の数 k が偶数の時 :

$$Ft(X) = X_1 / \sqrt{2} + X_2 \sin(t) + X_3 \cos t + X_4 \sin(2t) + X_5 \cos(2t) + \dots + X_k \sin(kt/2)$$

・変数の数 k が奇数の時 :

$$Ft(X) = X_1 / \sqrt{2} + X_2 \sin(t) + X_3 \cos t + X_4 \sin(2t) + X_5 \cos(2t) + \dots + X_k \sin((k-1)t/2)$$

ここで t の値は $-\pi < t < \sim +\pi$ の範囲で変化する。

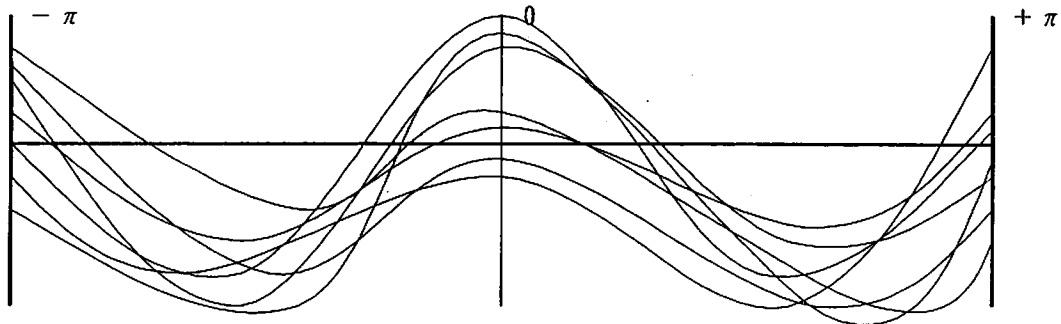


図 5 - 3 三角多項式グラフ

図では曲線の本一本が一個のパターンに相当する。このグラフはパターンが時系列データに置き換えられている点で他のディスプレイ手法とは際立って異なっている。

参考文献)

Andrews, D. F., "Plots of high-dimensional data", Biometrics, 28, 125-136 (1972).

5. 2. 2 計算機のグラフィック表示機能を生かした表示手法

現在ではパターン認識/多変量解析/統計解析に計算機は様々な場面で不可欠なものとなっている。単に計算するのみならず、計算結果の速やかな表示も計算機のグラフィック技術の発展とともに急速に展開されつつある。先に述べた一連の手法は総て人間がチャートを見て、最終的な判断を下す作業を助けるという事を目的とする手法である。結果の図形出力は重要であり、この点からも計算機のグラフィック機能との連携は極めて好ましいものである。

このような事実を背景として、計算結果のグラフィック表示が注目されつつある。特にグラフィック表示機能は日進月歩であり、計算機を単なる表示手段として考えるのではなく、さらに一歩進め、計算機中心或いは計算機特有の機能を積極的に活用した表示手法を展開する事が今後の表示手法に大きな影響を及ぼすものと考えられる。

残念ながら現時点では計算機の表示機能を積極的に扱った表示手法は少ない。単に表示画面を紙からグラフィックに置き換えただけのものにしかすぎない。大部分は従来の解析手法の結果の表示という補助的利用に止まっている。ここでは計算機のグラフィック機能を積極的に利用する表示手法はこうあるべきであるという原則論と、著者が試みている表示手法を例にとり簡単に説明する。

まず最初に、グラフィックを積極的に用いた表示の特徴を以下にまとめる。

- ① 結果のグラフィック表示が簡単になる
- ② 従来は困難であった様々な形式での表示が可能となる。
 - ・結果の3次元表示 (立体視による/3次元ディスプレイによる)
 - 従来は紙の上でしか表示出来なかった為に、表示出来るデータは物理的限界から

2次元に限定されていた。3次元プロットも可能ではあったが、3次元データを2次元上で眺めるために、立体的な奥行きをつかみにくかった。

- ・ 計算機でしか出来ない表示も可能となる。
画面の動的移動・変化／重ねあわせ／透視変換等の機能を利用した表示
- ・ 計算機の諸機能の利用
カラー表示機能／半透明機能等を駆使した表示

□ 回転チャートによる計算機上でのパターン表示及び1次元への次元圧縮

計算機特有の表示手法として幾つか考えられる。ここでは著者が考案した回転チャートについて述べる。この回転チャートは数値データの数に応じた複数の扇形より構成される円盤である。また、数値データの値は色と対応されており、個々の扇型は対応する数値データ固有の色にぬられている。パターン間の比較はこの円盤を比較する事で容易に実行可能である（静止画による比較）。この状態での比較は従来手法の表示手法と大きな差異はない。

この円盤をグラフィックディスプレイ上で回転させる事で円盤はそのパターンに固有の色となる（動画による動的比較）。この円盤を回転させる事がこの回転チャートの最大の特徴である。回転した状態の円盤の色を比較する事で多次元情報が簡単に1次元（色）情報に変換されたことになる。この結果、パターン相互の比較を簡単且つ正確に行う事が可能となる。これは、この回転チャートが本章で紹介した他の表示手法には存在しない、次元圧縮機能（N次元→1次元）を持っている事を意味する。このようにグラフィックディスプレイの特徴を利用する事で、従来の静的な表示手法にはない特徴を備えた表示手法を実行する事が可能となる。

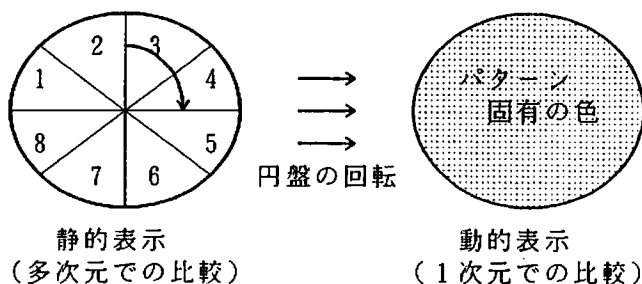


図5-4 回転チャートによるパターン表示及び次元圧縮

図4は8次元パターンの回転チャートを示している。円の内部の8個の扇形はパラメータの値の大きさに対応して色付けられている。このチャートはこのままでもパターン間の比較は可能であるが、パターン内の複数情報を意識しながら互いに比較する事が必要な為、視認性に難点がある（これは多次元での比較を行うためであり、他の表示手法にも共通に存在する）。つぎに、このチャートを計算機上で回転させる事で回転チャートはそのパターン固有の色となって見える。この色は用いた複数データの総合としての色であり、利用者はこの色の変化を認識するだけでパターン間の差異を認識する事が可能である。即ち、複数の色（多次元）で構成される情報が一色の1次元情報に変化した事を意味し、1次元での比較は簡単であり視認性も向上する。

□ グラフィックディスプレイの機能を利用した

立体レーダーチャート（テントチャート）

グラフィックディスプレイ上では回転／移動等を行う事が容易である。この機能を利用することで、従来から利用されてきた表示手法を改善する事ができる。

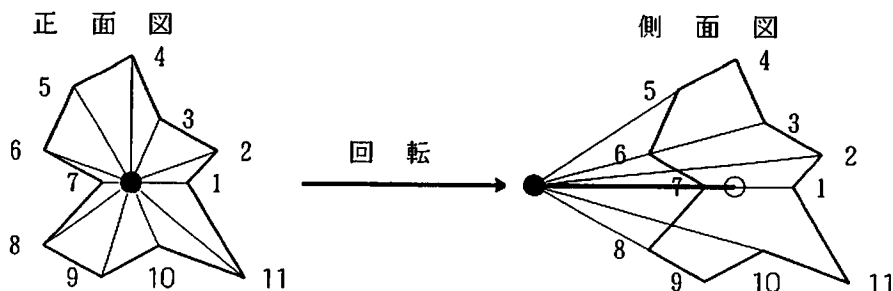


図5-5 レーダーチャートに奥行きを与えた3次元立体レーダーチャート

例えば、2次元のレーダーチャートに奥行きを与え、3次元にした立体レーダーチャート（テントチャート）がある。この3次元立体レーダーチャートは3次元目となる奥行きに目的変数をとる事で、目的変数と従属変数とを同時に表示する事が可能である（図5-5）。このチャートを用いれば、ディスプレイ上で回転する事で簡単に目的変数と従属変数との比較ができる。

5. 4. パターン空間におけるパターンの相互関係を図示する為の手法

前章で述べた表示手法は個々のパターン自体の情報を表示する手法であり、パターンを識別することを目的とした手法である。もちろん、パターン同士の相互関係を認識することも可能であるが、このパターン間の相互関係は人間が表示図形をみる事で認識する事が基本となっており、明確な形で捕らえるという事は困難である。

このような表示手法に対し、パターン間の相互関係（パターン空間上における相互位置関係）を明確にする事を目的とする手法が存在する。ここではこのようなパターン間の相互関係を重視して表示する手法について述べる。

5. 4. 1 パターン間の相互関係の表示手法について

パターン間の相互関係を重視して表示するアプローチはおおきく2種類に分類される。一つは投影（PROJECTION）と呼ばれ、もう一つはマッピング（MAPPING）と呼ばれる。以下、おのおのについて簡単に説明する。

□ プロジェクション (PROJECTION)

プロジェクション（投影）は一般的にはリニアプロジェクションとよばれる。このリニアプロジェクションはパターン空間中におけるパターンの存在状態には手を触れず（強制的に新しいパターン空間に変換する事はしない）、単にパターン空間を透視する方向（なるべく分散の大きな方向からの透視を目指す）について検討を加えるという手法である。

このアプローチには、先に述べた主成分分析より得られた結果を2次元散布図として表示（BI PLOT図）する手法がよく使われている。これは、主成分分析の特性により第1及び第2主成分をそれぞれX、Y軸として散布図を作成すれば、多次元空間におけるパターンの相互関係を2次元だけではあるが、最も広がりのある方向からの透視が保証される為である。

□ マッピング (MAPPING)

マッピング（写像）は多次元空間におけるパターン相互の位置関係を保ちつつ、人間が認識しうる2/3次元空間における位置関係へと変換することを意味する。従って、前記リニアプロジェクションと異なり、多次元空間中の個々のパターンについて新たに2次元のパターンベクトル（2/3次元）を算出することが必要である。

ここではこのようなマッピング手法の中で、最も典型的なもの「ノンリニアマッピング（NLM: NON-LINEAR MAPPING）」を紹介する。この手法はサモン（Sammon, J. S.）により提唱された手法である。

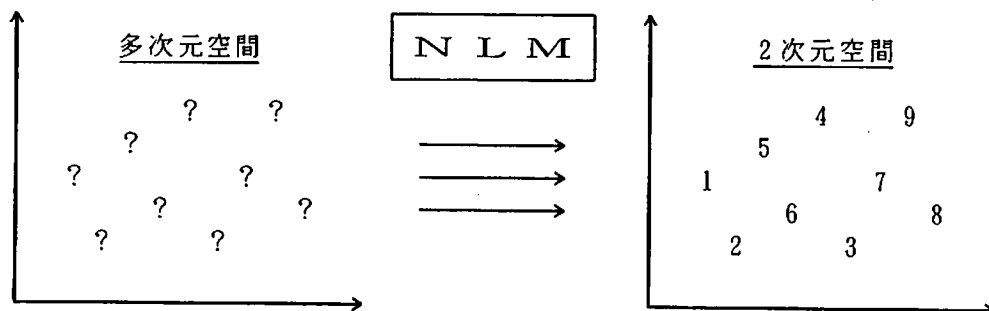


図5-6 ノンリニアマッピングによる多次元空間の2次元空間への写像

ノンリニアマッピング (NLM) 計算手法

NLMでは多次元空間上のパターンベクトル（PA、PB）を、ある一定の拘束条件の下で2次元空間上のパターンベクトル（P2A、P2B）に変換する事が必要である。

多次元 (d次元) 空間上のパターン A, B

$$\begin{aligned} PA &= (X_{A1}, X_{A2}, X_{A3}, \dots, X_{A, d-1}, X_{A, d}) \\ PB &= (X_{B1}, X_{B2}, X_{B3}, \dots, X_{B, d-1}, X_{B, d}) \end{aligned}$$

2次元空間上のパターン M, N

$$\begin{aligned} P2A &= (Y_{A1}, Y_{A2}) \\ P2B &= (Y_{B1}, Y_{B2}) \end{aligned}$$

パターン M と N の多次元上における距離関係を 2次元上に写像 (Mapping) する為には、多次元空間上のパターン間の距離 DX_{AB} と 2次元空間上のパターン間の距離 $D2_{AB}$ との値が同じになる事が必要である。

$$DX_{AB} = \sum_{i=1}^d \{(X_{Ai} - X_{Bi})\}^{1/2} \quad ()$$

$$D2_{AB} = \sum_{i=1}^2 \{(Y_{Ai} - Y_{Bi})\}^{1/2} \quad ()$$

ここで、 DX_{AB} は多次元空間上におけるパターン A と B 間の距離を、 $D2_{AB}$ は 2次元空間上におけるパターン A と B 間の距離を表す。多次元空間上での関係を 2次元空間に写像する事は、 DX_{AB} と $D2_{AB}$ が同じ値となることを意味する。従って、 DX_{AB} と $D2_{AB}$ とを用いてエラー関数 E_i を設ける。このエラー関数 E_i が 0 になった時、多次元空間の距離関係が 2次元空間上に完全に写像された事になる。

$$E_i = \frac{1}{\sum_{A \neq B} DX_{AB}} \sum_{A \neq B} \frac{n (DX_{AB} - D2_{AB})^2}{DX_{AB}} \quad (i = 1, 2, \dots, n)$$

ここで、n はパターンの数を示す。

このエラー関数の極小化は (Y_{i1}, Y_{i2}) ($i = 1 \sim n$) の初期チャートを適当に与え

$$\frac{\partial E(Y_{ij})}{\partial Y_{ij}} = 0 \quad (i = 1 \sim n; \quad j = 1, 2) \quad ()$$

の値をニュートン・ラプソン (Newton-Raphson) 法、或いは最急降下法等の手法を用いて最適解を求めれば良い。初期チャートとしてはどのようなプロット図でも利用可能である。しかし、最適化に要する時間や最小/極小問題を考えるならば、全パターンが 2次元上に適当に分散していることが理想である。このような観点から、NLM の初期チャートには主成分分析で得られた第 1 及び第 2 主成分による散布図を用いるのが一般的である。

この他にサモンは以下に示されるような簡単な漸化式を用いている。

$$Y_{ij}(m+1) = Y_{ij}(m) - k \Delta_{ij}(m) \quad ()$$

$$\text{ここで } \Delta_{ij}(m) = \frac{\partial E(m)}{\partial Y_{ij}(m)} \left/ \left| \frac{\partial^2 E(m)}{\partial Y_{ij}(m)} \right|^2 \right. \quad \text{で、} k \approx 0.3 \sim 0.4$$

また $Y_{ij}(m)$ 、 $E(m)$ は m 回反復計算を行った後の値である。

理論的にはこのエラー関数を 0 にする事は不可能であるが、この値を 0 に近づける事は可能である。従ってこのエラー関数の値が最小となるように各パターンの 2次元座標を決定することで写像がおこなわれる。

参考文献)

Sammon, J.S., "A non-linear mapping for data structure analysis", IEEE Transaction on Computers, C 18, 401-409 (1969).

ノンリニアマッピングにおける留意事項

ノンリニアマッピングによる計算において、エラー関数の値が0となる事は理論上あり得ない。従って、解析はエラー関数の値が解析に支障のない程度に小さくなった時のマップ(散布図)を用いる事になる。この値に関する基準はない。また、このエラー関数の最適化過程では、この種の最適化問題では常に問題となっているローカルミニマの問題がある事も見逃せない。現在のマップが最良のものとは限らず、むしろローカルミニマに落ち込んでいる可能性の方が高い事を解析時には常に意識すべきである。

5.4.2 ニューラルネットワーク(バックプロパゲーション)による写像

写像を行うアプローチとしてニューラルネットワークを利用することが可能である。この写像はニューラルネットワークのバックプロパゲーション手法を利用するもので、バックプロパゲーションの中間層からの出力値を写像されたデータとして用いる。

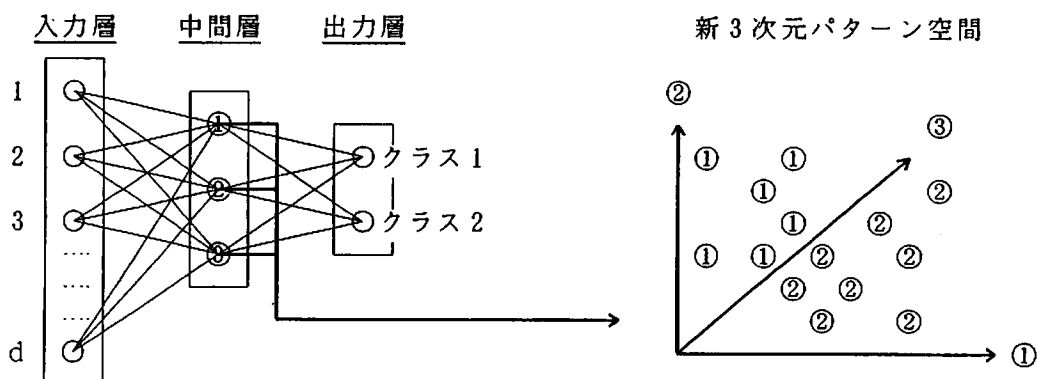


図5-7. バックプロパゲーションによる写像。 d次元⇒3次元
(新たな次元は各パターンがクラス1とクラス2とに分類されやすいように3次元空間上に分布していることに注意)

バックプロパゲーションによる写像は他のアプローチで得られる写像と異なり幾つかの点立った特徴を持つ。

- (1) 写像先の次元を自由に設定できる。(次元拡大、等次元写像(投影)、次元圧縮)
- (2) 写像される時、教師データの情報を加味した写像が行われる。

これらの特徴について順に説明する。

バックプロパゲーションによる次元拡大、次元圧縮

バックプロパゲーションによる写像では写像先の次元を自由に設定することが可能である。この設定はバックプロパゲーションの中間層のユニット数を指定することで簡単に実現される。従って初期の次元よりも大きくも小さくも自由にコントロール出来、解析目的に従って柔軟に対応出来る。

次元圧縮という目的では2次元に限定されたノンリニアマッピングと比べ、自由に次元を設定できるという点でより汎用性の高い写像手法となる。さらに、主成分分析で行ったのと同様に次元減少という目的にも利用出来る。

教師データの情報に従った写像

本アプローチによる次元拡大/圧縮において、新たに形成される次元はバックプロパゲーションの学習過程で得られた情報(教師データの情報)を反映している。例えば、バックプロパゲーションでクラスデータが教師データとして利用されたならば、そのバックプロパゲーションから得られる写像データは個々のパターンを教師データに従って分類するのに都合の良いようなパターン分布を形成したデータとなる。この点が本アプローチによる最大の特徴である。

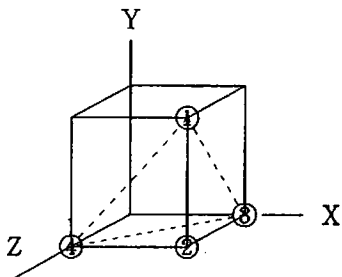
従来手法によるアプローチでは元のデータが持つ情報を失わないということの特徴としてきた。主成分分析ではd次元空間中の分布状態を視点を変えてみているだけで、d次元空間を変形させることはしない。また、ノンリニアマッピングはd次元を2次元に落とす時、d次元空間におけるパターン間の相互位置関係を保持するようにしている。このようなアプローチにたいし、バックプロパゲーションによる次元拡大縮では元のデータを教師データに従って積極的に変化させる。この特徴を利用することで、従来手法では得られなかった様々な解析を実現することが可能となるだろう。

5. 5. 1次元への次元圧縮による写像

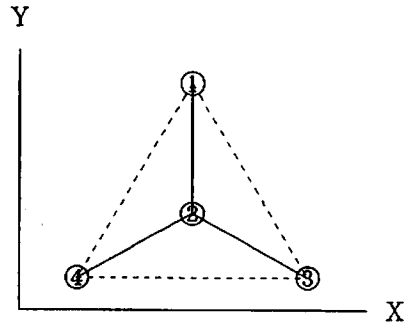
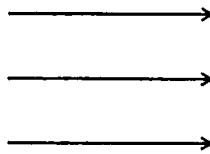
マッピングでは主として2次元での表示が中心である。この2次元はパターン間の複雑な相互関係を紙上で表現しうる、また人間が簡単に認識できる最大の次元であるためである。更に1次元への圧縮が出来ればパターン間の相互関係という貴重な情報は失われるが、パターン同士の絶対比較という点では便利な機能である。1次元への次元圧縮過程で多くの情報が失われるという点を見捨てるならばこの問題は数学的には簡単な問題である。単に数値データに置き換えるだけでなく、例えば前節で述べた計算機による回転チャート等、より適切な表示手法等も考案する事で今後の展開が期待される。

QUIZ

設問： 3次元上の点①～④を2次元上にマッピングせよ。



- ① (1, 1, 1)
- ② (1, 0, 1)
- ③ (1, 0, 0)
- ④ (0, 0, 1)



- ① (?、?)
- ② (?、?)
- ③ (?、?)
- ④ (?、?)

* 写像の基準点は②を原点にせよ

5. 多次元空間を見直すためのアプローチ

総てのパターン認識の最終目的は、多次元空間上に分散しているパターンの状態を解析することである。パターン空間上ですべてのパターンをクラス毎に分類できるような面を探すのがクラス分類手法であり、このアプローチについてはすでに説明した。

パターン認識のアプローチとして、このパターン空間を異なる角度から見直すという試みも存在する。パターン空間を見直すことで従来は見えなかったものが見えるようになり、不可能だったことが可能となるといったことが期待できる。ここではこのような、パターン空間を見直すことを主眼としたアプローチに付いてまとめてみる。

5.1 主成分分析（視点を変えてみるアプローチ）の基本概念

多次元上に存在する個々のデータ間の相互距離関係を新たな視点（座標系）より眺め直すことで、データ間に存在するなんらかの分散要因を探し出す手法である。従って、主成分分析は後にのべるマッピング手法の一つとみなされる。

主成分分析において新たに形成される座標系はパターン空間中におけるパターンの広がりを最も大きな広がりをもつ次元から順番に取り出される。主成分分析とは多次元空間上の全パターンの分散を最大に見せる新座標系を求める手法である。

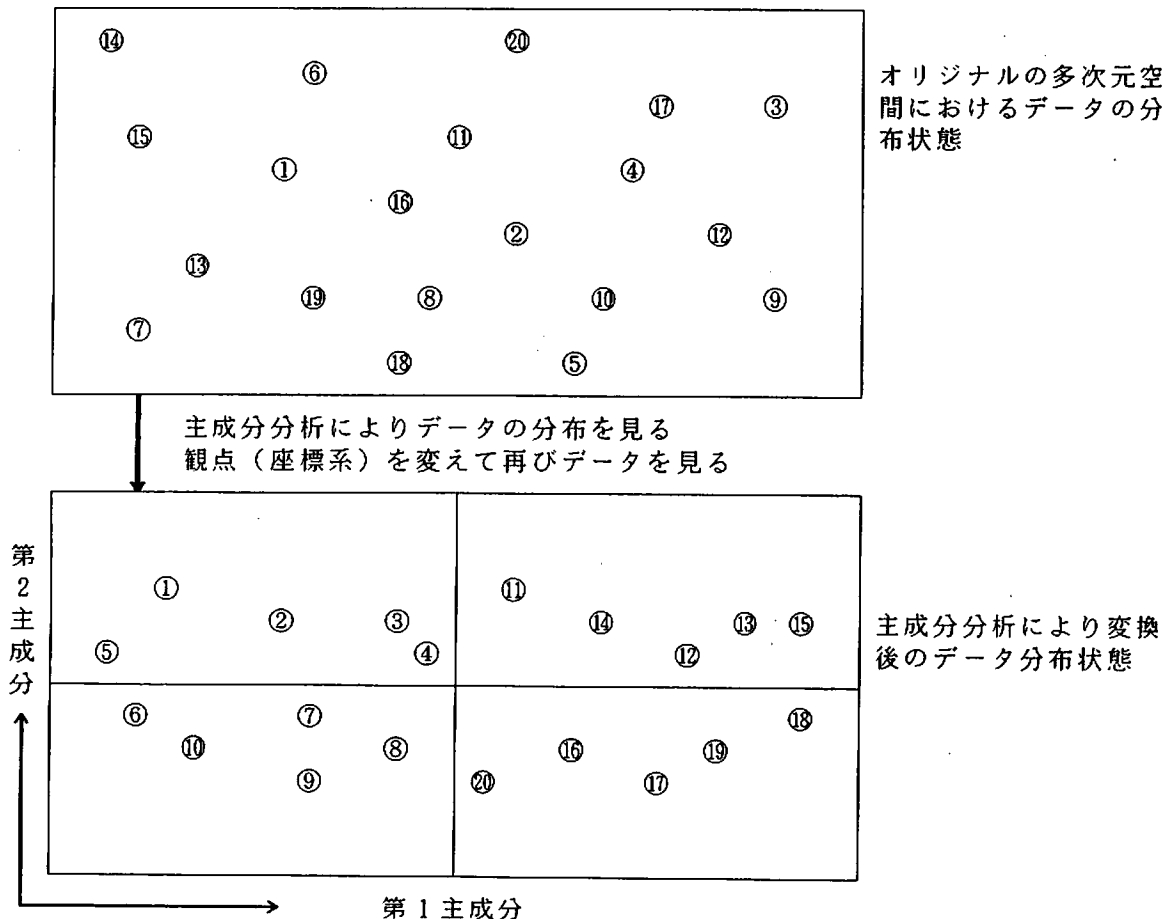
*1. 単にパターン分布の見方を変えるだけなので、オリジナルパターンの次元数と解析結果の次元数（主成分軸の数）とに変化はない。

*2. 通常人間は物体を正確に認識しようとする時には上下左右に視点を変え、物体を確実に認識できるようにする。この考えは主成分分析でも同じである。

一枚の紙を例にとって考える。紙とわかっている時、斜めからみてもその特質は推測可能である。しかし、紙と分かっていない時は困難である。理想的には、正面図（最大平面：XY平面）と側面図（直線：Z軸）をみる事で紙の本質がつかめるであろう。主成分分析は3次元以上の多次元空間上におけるパターンの広がりを最も正確に認識出来る（パターンの分散が最大となる）新たな座標系を見出す手法である。

この主成分分析による2次元散布図に基づき、個々のデータの分布状態を検討する。

この時、X軸（第1主成分）及びY軸（第2主成分）について、それぞれのデータがどのような傾向で分布しているかを検討する。この傾向（分散要因）を取り出す事が主成分分析の最終目標である。



X軸（第1主成分）： ①～⑩と⑪～⑳とを分離するのに重要な因子である。
 Y軸（第2主成分）： (①～⑤、⑪～⑮)と(⑥～⑩、⑯～⑳)とを分離するのに重要な因子である。

図1. 主成分分析概念図

主成分分析出力結果解析例)

①～⑩はグラム陽性菌に活性ある化合物、⑪～⑳はグラム陰性菌に活性ある化合物である時、

結論 ——— X軸（第1主成分）はグラム陽/陰性菌に関係する次元である。

①～⑤と⑪～⑮は酸性抗菌剤であり、⑥～⑩と⑯～⑳はその他の抗菌剤である。

結論 ——— Y軸（第2主成分）は化合物の酸性等の物性に関する次元である。

□ 寄与率及び累積寄与率について留意点

主成分分析で求められる寄与率とは、個々の主成分軸（新次元軸）が有する分散量と、オリジナルデータの全分散量との比率をいう。また、累積寄与率とは複数の主成分軸の寄与率を加えた値である。全主成分軸について寄与率を合わせれば100%となる。

寄与率の高い主成分軸はd次元パターン空間上におけるパターンの分散を効率よく説明し、パターンの分散の理由等を発見する要因解析に対し非常に貴重な情報を持つ次元である。反対に寄与率の低い主成分軸は情報量の少ない軸となる。

例えば、5個の記述子で示される5次元データを用いて主成分分析を行った時、以下に示す寄与率が得られたとする。

第1主成分 = 100%、 第2、3、4、5主成分 = 0%

この結果は、用いたデータセットが第1主成分のみで総てのパターン分布の広がりを見しうる事を意味している。即ち、全パターンは一直線上に分散して存在していることになる。ちなみに他の主成分軸は寄与率が0%なので、総てのパターンがただ一点に集中している。

以上の事実より、多次元データを2次元の散布図を用いてデータの広がりを観察するならば、第1及び第2主成分を合わせた累積寄与率が高いほど正確にデータの分布状態を把握することが可能となる。

□ 記述子の因子負荷量

因子負荷量とは各記述子が個々の主成分軸に及ぼす影響の程度を表す数値である。この因子負荷量について図2のバイプロット(BIPLLOT)図を例に取り、各記述子(数値データ)の第1、第2主成分への寄与の程度(因子負荷量)を以下に示す。

記述子1 記述子2 記述子3 記述子4 記述子5

第1主成分	中	中	大	小	大
第2主成分	大	小	中	小	中

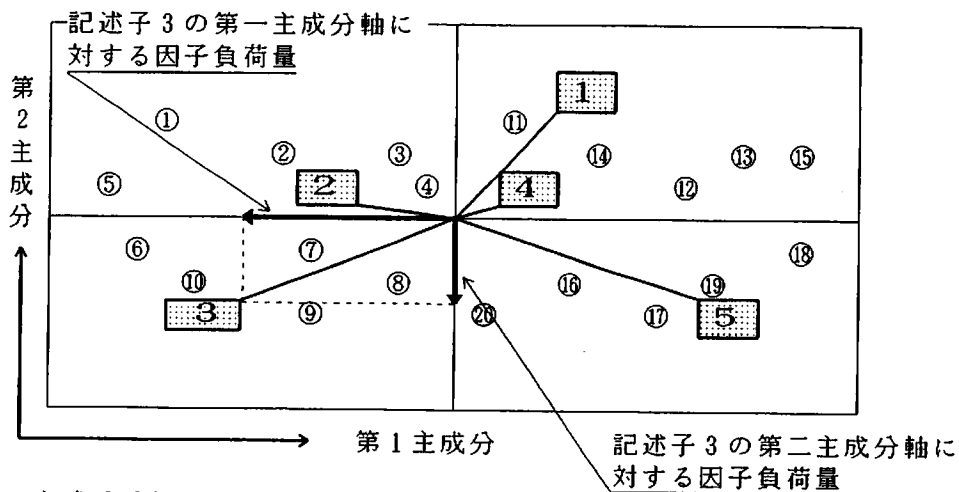


図2. 主成分分析結果のバイプロット図

バイプロット図では解析に用いた記述子の各主成分軸に対する因子負荷量が2次元散布図上に示される。中心から線が引かれ、四角の箱が示されているが、この線の長さや方向がX軸（第1主成分）とY軸（第2主成分）の因子負荷量に関する情報をベクトル表示したものである。従って、この線（ベクトル）の第1/2主成分軸に対する長さが、因子負荷量の各主成分軸に対する値を示している。尚、図中四角に囲まれた値は記述子の識別番号である。

因子負荷量は具体的にいうならば、用いた記述子と求められた主成分軸との相関係数である。従って、正/負のサインを持ち、絶対値は1を最大とする。この因子負荷量は記述子と主成分軸との相関を示す値であり、記述子を構成する1次元上の各パターンが新たに計算された主成分軸のどの位置にくるかを定める重要な要因となる。この因子負荷量が1の時、記述子上での各パターンの相対関係はそのまま保持された状態で主成分軸上に転写される。値が-1の時、この関係は正反対となる。値が0の時、記述子と主成分軸との間には何の相関も存在しない事になる。

5. 2 主成分分析法の様々な利用形態

□ 寄与率による次元減少

主成分分析法を利用する事で次元減少を実現する事が可能である。これは主成分分析で得られた個々の主成分に関する主成分得点を新たな次元データとする事で実現される。この時、総ての主成分を用いずに累積寄与率が十分に高くなり、これ以上主成分を増やしても累積寄与率に大きな変化が認められなくなった時点までの主成分を用いることで、情報量の減少を伴う事なく次元減少をおこなう事ができる。

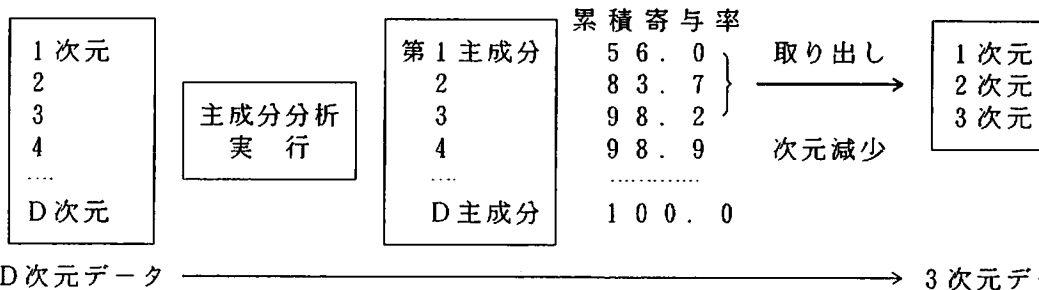


図3. 主成分分析によるD次元データから3次元データへの次元減少
全体がD次元のデータで表現されるパターン空間を、主成分分析により分散全体の98.2%を保持しつつ3次元パターン空間へと変換する。

主成分分析を用いた次元減少はパターン全体の情報量（分散）を減少する事なく、簡単に次元を減少できるので実用価値の高いアプローチと言える。但し、この次元減少はノイズデータを含んだ形での次元変換及び減少である。一方、後に述べる特徴抽出手法はノイズを取り去る事による直接的次元減少である。表1に主成分分析による次元減少と特徴抽出による一般的な次元減少との差異をまとめる。

表1. 主成分分析による次元減少と特徴抽出による次元減少の差異

内容	主成分分析	特徴抽出
次元数減少	次元圧縮	次元カット
情報	情報圧縮	情報カット（ノイズカット）
情報内容	全く別のものとなる	選択された情報
情報量	変化せず（ノイズも含まれる）	減少する（ノイズがなくなる）
次元データ	新たに発生したデータ	特徴抽出で残ったデータ利用

以上の特徴を考慮するならば本アプローチは分類だけを目標とする時や、強力な特徴抽出手法が存在しない時には極めて貴重な次元減少手法となる。しかし、ノイズカットを重視し、情報内容の読み取りを重視するアプローチではその適用に留意が必要である。例えば構造-活性相関解析では、解析に用いた記述子（数値データ）の情報を取り出す事が最重要目的である。このような解析では、主成分分析で新たに発生した主成分の持つ情報の解釈が困難となる事が多く、効率的な構造-活性相関を行う事は困難である。

この次元圧縮の技術と線型重回帰手法とを組み合わせる事で、5章3節で述べるPCR法（PRINCIPAL COMPONENT REGRESSION）という新しい解析手法となる。

□ 因子負荷量による特徴抽出（記述子選択）

因子負荷量を見ることで特徴抽出を行うことが可能である。因子負荷量が各記述子のパターン分散に関与する割合を示すものであるならば、この割合が小さな記述子はパターンの分類等に重要な記述子とはいえない。

即ち、因子負荷量が小さい（バイプロット図で中心から伸びる線が小さなもの）記述子はノイズデータとして取り除くことが可能である。例えば図2中、記述子2は因子負荷量が第1/2主成分ともに小さく、ノイズとして最初に取り出されるべき記述子である。

また、このバイプロット図で各クラス毎のパターンの分布状態をみることで主成分軸が持つ情報を推察することができることは既に述べた。従って、解析目的に最も重要と思われる主成分軸に対し最も大きな因子負荷量を持つ記述子が最も重要な記述子であり、その他の主成分軸に対し大きな因子負荷量を持つ記述子があったとしても目的とする主成分軸に対する因子負荷量が小さければその記述子はノイズデータである。

5.3 主成分法的な考えを基本とした様々なアプローチ

現在、主成分法を基本としてその他の解析手法と結び付けて新たな解析手法とするアプローチや主成分法の基本的考えを若干修正して新たな解析手法に展開するアプローチ、あるいはその両方のアプローチをとるもの等があらたに展開されつつある。ここではこれらのアプローチについて簡単にのべる。

□ 他の手法との組み合わせによるアプローチ

PCR（主成分回帰）法（PRINCIPAL COMPONENT REGRESSION）

主成分分析手法と重回帰手法とを組み合わせた手法である。手続的には、最初に主成分分析による次元変換を行った後、その主成分を重回帰の説明変数として解析を行うものである。手法的には重回帰手法の一種と考えられる。

この解析作業を効率的に行うために、重回帰の説明変数としてどの主成分をもちいるかの判断が必要となる。一般的には最初の数個の主成分で、十分に大きな累積寄与率となる数の主成分を用いている。

PCR（主成分回帰）法の利点及び欠点

本手法をもちいる事による利点は大きく2つ考えられる。

①パターンよりも次元数が大きいデータ（ケモメトリクス分野での解析に多い）を用いた線型重回帰の解析が可能となる。

パターン認識を行うにあたり常に留意すべきこととして、解析結果の信頼性を高く保つという事がある。この解析結果の信頼性を高く保つ条件として現在いくつかの制限条件（詳細は第4章にて説明する）が知られている。例えば、線形重回帰を行う時にはサンプル数が解析に用いた記述子数の4倍以上なければならないといったものがある。このような条件を満たすことが出来ないデータを用いて解析することは出来ない。このようなデータを用いる時、最初に主成分分析を用いて次元数をへらすことで、先の制限条件を満たして解析することが可能となる。

②共線性（多重相関）の高いデータに対しても良好な結果が得られる。

使用データ中に共線性の高いデータがあったとしても、一旦主成分分析により変換されたデータを用いて解析を行うならば共線性に関する問題を回避することが可能である。具体的には寄与率の高い主成分を主体としたデータを用いて線型重回帰を行う事となる。

一般的に、機器スペクトルから得られるデータはパターン数よりも次元数の方が大きいことが多く、また一部スペクトル（近赤外、UV等）では共線性の高いデータとなることが多いので、本手法はケモメトリクス分野での適用が有効である。

PCR手法における欠点として以下の事実がある。

①最終回帰式の係数の比較や各変数を持つ情報の解釈が困難である。

線型重回帰で用いる数値データとして主成分分析で変換されたデータを用いる（各変数は最初に用いた数値データの線型結合で示されている）ため、最終回帰式の係数だけから最初に用いた変数に関する情報を取り出す事は不可能である。一旦、元のデータを用いた回帰式に変換した上で回帰係数に関して解析する事が必要となる。

このPCR法を行う時の解析作業の流れ図を図4に示す。

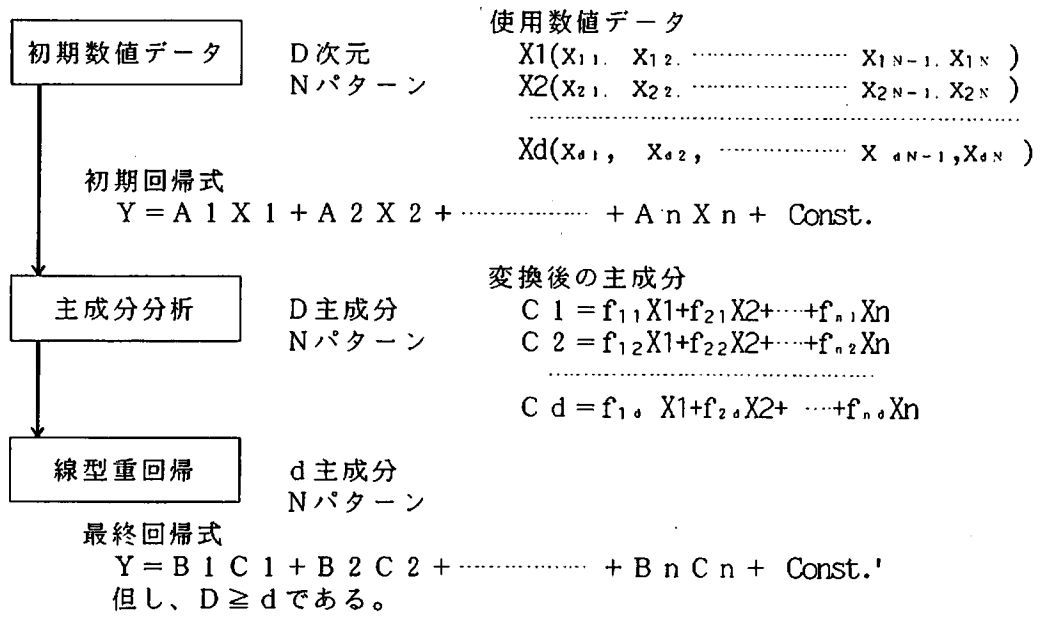


図4. PCR法の解析作業流れ図

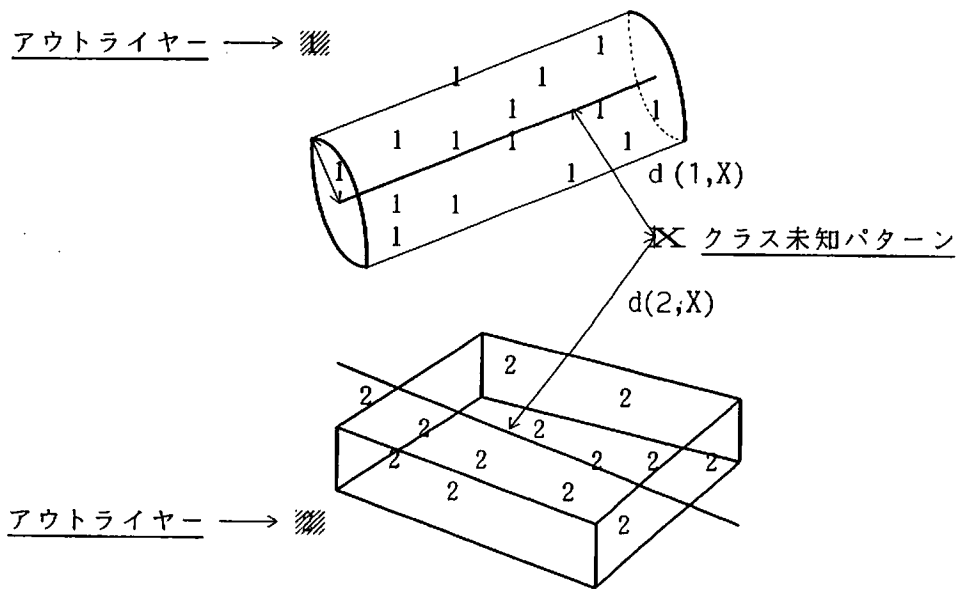
- 主成分法を基本とし、重回帰的アプローチを可能としたアプローチ
SIMCA法 (Statistical Isolinear Multiple Components Analysis)
(Soft Independent Modelling of Class Analogy)
(SIMple Classification Algorithm)
(SIMilarographic Computer Analysis)

S. Woldが提唱したもので、主成分的考えを利用した分類手法である。主成分分析はクラス情報を用いない教師無し (UN-SUPERVISED) アプローチであるが、SIMCA法ではクラス情報をもちいる教師有り (SUPERVISED) アプローチを取る。即ち、SIMCA法の解析ではクラス情報を参照することで、クラス単位でのパターン分布状態を吟味する。この結果、主成分分析では出来なかったクラス分類 (多クラス分類可能) が可能となり、さらには重回帰法的解析も可能となる。

SIMCA法ではクラス毎に主成分分析法が適用されてクラスモデルが作成される。

$$Y_{ik}^{(q)} = \alpha_i^{(q)} + \sum_{a=1}^d \beta_{ia}^{(q)} \theta_{ak}^{(q)} + \epsilon_{ik}^{(q)} \quad ()$$

Y_{ik}はパターンiのk番目のパラメータを示す。
ここでqはパターンがクラスqに属する事を示す。
以下には、SIMCA法にて計算されたクラス1と2のパターン群のイメージ図を示す



SIMCA法では単にパターンの分布状態に関する情報を提供するのみならず、様々な情報を提供する。例えば解析に用いられる記述子に関する情報は特徴抽出等の手法に利用される。また、クラス分類の為の情報等である。

アウトライヤーの概念

アウトライヤーの概念はSIMCA法では、前図で示した超円筒及び超長方体の中に該当するパターンが存在するか否かにより決定される。この構造体の中にパターンが存在しない時、そのパターンはアウトライヤーとみなされる。

このアウトライヤーを決めるに必要な構造体の大きさは、分類時にある程度人間が決定する事が可能である。

アウトライヤーの分類時の扱い

このようなパターンは分類/予測等においては解析対象として扱われない(分類不可能)事になる。但し、このアウトライヤーの定義を無視して分類が行われる事もある。

モデリングパワー (Modelling Power)

この指標は解析に用いられたデータの全分散量のうち、該当する数値データ(記述子)によりどの程度説明されているかを示すものである。従ってこの指標により、個々の記述子が持っている分散に対する貢献分の程度かわかる事になる。この値は以下の式により求められ、その値が1に近い時モデリングパワーは高く、0に近くなるにつれ低くなるものである。

$$\Psi_i = 1 - (S_i / S_{i,y}) \quad ()$$

i は用いる数値データ(記述子)のID番号ですべてのクラスに関して計算される。

S_i は記述子 i の総てのクラスに対する残差標準偏差(RSD)であり、 $S_{i,y}$ は母集団データ Y の標準偏差である。

モデリングパワーの値が0に近いものは、ノイズデータとして取り除く事で特徴抽出を行う事が可能となる。

ディスクリミナトリーパワー (Discriminatory Power)

この指標は個々の記述子単位で、各クラスのパターン同志がどの程度分離性よく離れているかを示す指標である。

いまパターンがクラス q と r とに別れているとすると、記述子 i のディスクリミナトリーパワーは以下の式でしめされる。

$$\Phi^{(i)} = \left[\frac{(q)^2}{S_{i,r}} + \frac{(r)^2}{S_{i,q}} \right]^{1/2} - 1 \quad ()$$

この値は0に近い程低いディスクリミナトリーパワーとなっている。従って、このような記述子はパターンの分離という観点ではノイズデータとなる。

クラス未知パターンの分類

クラス未知パターンの分類は図にて示されるようにクラス未知パターン X から各クラスの主成分軸に下ろした垂直2等分線の長さにより判定される。

即ち、主成分モデル面からの許容領域を示す残差標準偏差(RSD)を求める。このRSDと各パターンからクラスモデル面までの距離を計算して大小を比較する。この時F検定によりクラス判定の検定を行う。

また、各クラス毎に構造体が決定されている時はこの構造体にパターンが落ち込んでいいる時、その落ち込んだ構造体が帰属しているクラスにアサインされる。

□ PLS法 (PARTIAL LEAST SQUARES METHOD)

SIMCA法と同様、主成分分析法の改良アプローチである。

6. 最適化手法について

6. 1. 最適化手法一般概念

最適化は様々な分野で応用される基本的手法の一つである。例えばパターン認識では判別分析における判別関数の決定に利用される。化学関連分野では、最近注目を浴びている分子力学 (MOLECULAR MECHANICS) の分野で歪みエネルギーの最適化過程に利用されており、その他様々な場面でも利用されている。

この最適化に関し、幾つかのアプローチがある。ここではシンプレックス法とダイナミックプログラミングの2法について簡単にのべる。両手法とも従来行われてきた微分法 (最小二乗法) の欠点を補うものとして提唱されてきたものである。これら2手法について論じる前に従来手法 (偏微分法) による最適化の問題点を簡単にまとめてみる。

従来手法の問題点:

- ① 関数が微分不可能な時がある (領域の定められた最適化問題)
- ② 偏微分法の接線の傾きが0の状態は最小/極小状態だけでなく、変曲点でも発生し、最適状態の探索が困難になる
- ③ 評価関数の偏微分より得られる連立方程式の解決は困難な時がある
- ④ 制限条件がある時の最適化は困難である

このように従来のアプローチは幾つかの欠点を内包するものであったが、シンプレックス法やダイナミックプログラミングを用いる事でこれらの問題点を解決する事が可能となる。ここでは最小二乗法による最適化手法については述べない。関心のある方は、他の専門書を参照されたい。

6. 2. シンプレックス法 (SIMPLEX OPTIMIZATION) による最適化

□ シンプレックス法概要

シンプレックス最適化は、応答表面上 (評価関数にて形成される) に存在するシンプレックス (この場合トライアングル: 3角形) を一定のルールに従い反転させ、この反転を繰り返すことで最適値を探すものである。最適値を探すのに用いられる応答表面は、ある最適値をピークとする等高線上の起伏を持つ空間であるとする。この応答表面は、探索対象となる事象と対応がとれていることが必要である。

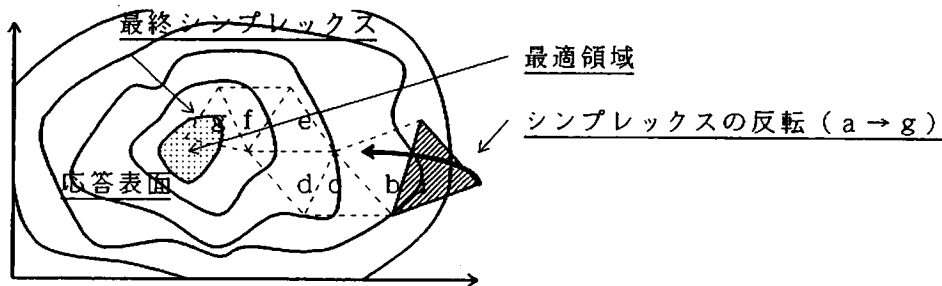


図6-1 シンプレックス法概念図

□ シンプレックス法による最適値探索ルール (反転ルール) について

シンプレックス法は応答表面上のシンプレックスを反転させて移動させる事で最適値を探索する手法である。従って、シンプレックスの反転方法がこの最適化手法において最も重要な過程となる。このシンプレックスの反転方法について、図2を用いて簡単に説明する。

(1) 初期シンプレックスの形成及び最初の反転による第2シンプレックスの形成

応答表面上に任意の3点 (W_1 , W_2 , W_3) を取り (第1シンプレックス)、この3点における評価関数の値を調べる。この時、評価関数の値が $W_1 < W_2 < W_3$ となったとする。この評価関数の値は小さい程望ましい値であるとした時、このシンプレックスは3点のうち望ましい値を示した上位2点の W_1 と W_2 を軸とし、最も望ましくなかった

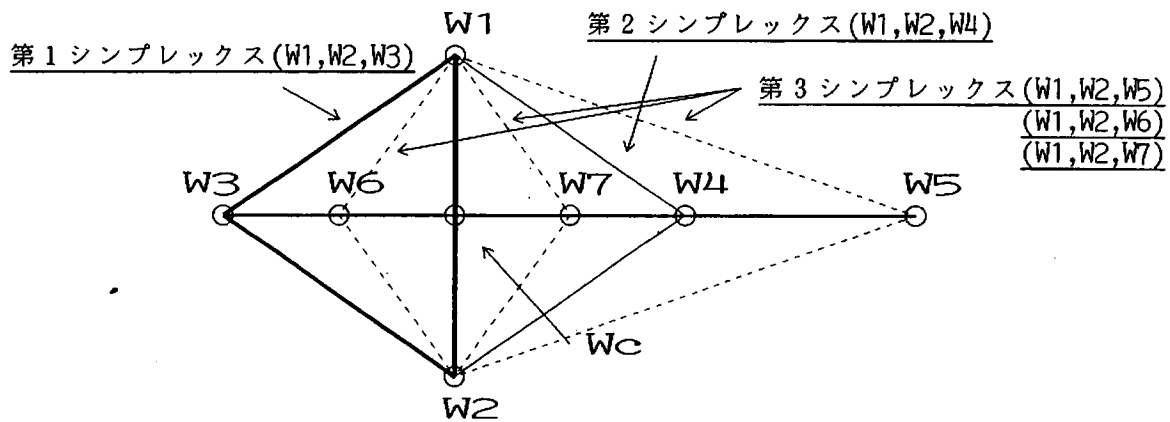


図6-2 シンプレックス反転に関するルール

W3を頂点としてWcと対称の方向に向かって反転させW4とする。この反転により、あらたなシンプレックス(W1、W2、W4) (第2シンプレックス)が形成される。ここで、

$$Wc = \frac{1}{d} \sum_{\substack{j=1 \\ j \neq 3}}^d W_j \quad (1)$$

$$W4 = Wc + (Wc - W3) \quad (2)$$

(2) 第2シンプレックスのW4について評価関数の値を調べる

この時の評価関数の値の結果に従って第3シンプレックスが形成される。この第3シンプレックスは、評価値が更に望ましい方向に変化した時と、望ましくない方向に変化した時とで異なったものが形成される。

① W4の評価関数が望ましい方向に変化した時 (W4が更に小さくなった時)。

あらたに作成されたW4の延長上に更に望ましい領域があるとし、その延長上に新たにシンプレックスを展開しW5 (第3シンプレックス) (W1, W2, W5)を形成する。

$$W5 = Wc + 2(Wc - W3) \quad (3)$$

② W4の評価関数が望ましくない方向に変化した時 (W4が大きくなった時)。

あらたに作成されたシンプレックスは望ましくないので、WcとW3との間にW6を設定し、新たなシンプレックス (第3シンプレックス) (W1, W2, W6)を形成する。

$$W6 = Wc - (Wc - W3) / 2 \quad (4)$$

③ W4の評価関数が、W3よりは良いが、W2より悪い時。

あらたに形成する頂点の位置はWcとW4との間に形成し、W7とする。新しいシンプレックスは(W1, W2, W7)より構成される。

$$W7 = Wc + (Wc - W3) / 2 \quad (5)$$

□ 初期シンプレックスの座標について

シンプレックス法は応答表面上のシンプレックスを移動させることで最適値を探すアプローチである。従って、この作業が成功するか否か、早く探せるのか否かといった問題は最初に置かれたシンプレックスの位置により大きく左右される事になる。この故に、シンプレックスの初期座標の決定は極めて重要なこととなる。現在迄に幾つかの初期座標の取り方が提案されている。この中で最もよく用いられている手法を以下に示す。

手順1. シンプレックスを決定する3点の内、最初の一点を決定する手法

シンプレックスの第一座標を決定する方法としては幾つかあるが、線型判別分析の判別関数の決定にシンプレックス手法を利用する時を例に取り、以下に代表的な3手法についてのべる。

- ①最初に各クラスの重心を求める。この重心間に結ばれた線とその線を2分割する垂直2等分線との交点を第1座標とする。
 - ②線型学習機械法等、他の解析手法で得られたウェイトベクトルを第1座標とする。
 - ③座標の総ての要素を+1か-1に設定したものをを用いる。
- ①、②の方法で第1シンプレックスを求める、最適値の発見を早める事が可能である。③の方法で行うと、手続きは簡単であるが最適値を見出すのに時間がかかる。

手順2. 第2、3番目の座標の決定

第1番目に求めた座標を示すウェイトベクトルの各要素にある一定の値(α)を加える事で決定する。この α の値を変える事で最適値を探すシンプレックスの大きさの指定が可能となる。

□ 第1シンプレックス作成事例

第1シンプレックスの座標を決定するベクトルが3次元の時を例とする。具体的な手順を図3に従って説明する。

最初に、初期シンプレックスを作成するのに必要な1個の頂点(C_1)を設定する。この設定は前節でのべたような手法で行われる。続いて、その第1座標の各次元(要素)にある一定値(α)を加える事で、新たに3個の頂点($C_2 \sim C_4$)が生成される。初期シンプレックスはこれら4個の頂点のうち3個を取り出して形成される。従ってこの場合、4種類の第1シンプレックスの形成が可能である。3次元以上のばあい、次元数+1の数だけ頂点が形成される。これらのうち3個を取り出して初期シンプレックスとすることができる。従って、より多数の第1シンプレックスを形成することが出来る。

- | | |
|--------------------------------|--------------------------------|
| 第1シンプレックス: (C_1, C_2, C_3) | 第3シンプレックス: (C_1, C_3, C_4) |
| 第2シンプレックス: (C_1, C_2, C_4) | 第4シンプレックス: (C_2, C_3, C_4) |

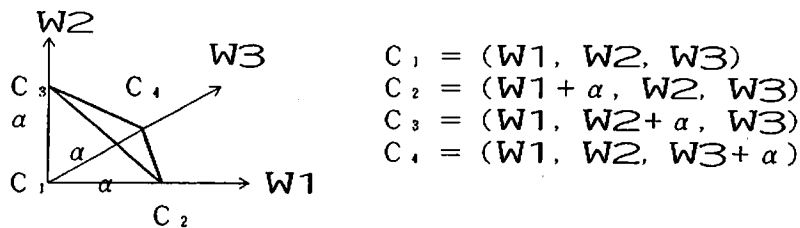


図6-3 初期シンプレックスの決定例

□ シンプレックス法における様々な特性

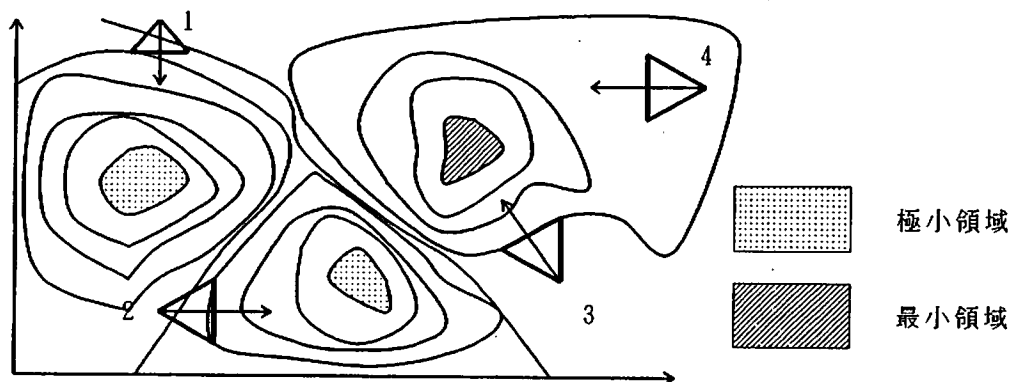


図6-4 シンプレックス手法におけるローカルミニマの問題

シンプレックス法を正しく使う上で幾つかの留意点があるので以下に明記する。

最小及び極小の問題

- (1) シンプレックスが最小点でなく極小点に陥る事がある
 応答表面上に唯1ヶ所の最小点だけが存在する時は問題ないが、通常は最小点の他に多数の極小点を有している。従って、通常のシンプレックスを動かす探索過程ではこの

極小点で探索が止まってしまうことが多い。現時点ではこの最小点を直接探し出す方法はない。従って、シプレックス法で探し出された領域が最小点か、極小点かを正しく見極める事が必要である。

例) 図4中1~4のシプレックス中最小点にたどり着く可能性のあるシプレックスは3と4のみである。1と2は極小点に行くことになる。

(2) 最小点、極小点の問題に関する対策

① 初期シプレックスの値を吟味する。

探索が最小点か極小点に行き着くかは、図からもわかるように初期シプレックスの位置に大きく影響される。この為、初期シプレックスの位置を吟味する事でこの問題点の改良が計られる。然し先にも述べたように、真の意味で最小点を探し出す手法は存在しないということは事実である。従って、最適値を探し出したとしてもその値が最小値なのか極小値なのかの確認/追試が必要である。

② 新たな最適値を探す

まず一つの最適値(極小値)を探し出し、この後に改めて次の最適値(極小/最小値)を探し出す。この手続きにはいくつかのアプローチがある。

i. 初期シプレックスを複数作成する。

最も簡単かつ効果的なアプローチは新たにシプレックスを作成しなおす事である。この時、シプレックスの値は大幅に異なっている事が必要である。さもなければ、再び同じ最適値にたどり着く事になりかねないので注意が必要である。

ii. シプレックスのサイズ変更による極小点からの脱出

シプレックスのサイズを大きくして、一旦たどりついた最小/極小領域を脱出させる。脱出に成功すれば、新たな最小/極小領域を発見出来る可能性が高い。この時、シプレックスのサイズは最小/極小領域を越えるのみならず、その領域を頂点とするまわりの等高線部分も凌駕したものでなければ、もとの最小/極小領域に戻ることになるので注意が必要である。

シプレックス法の探索速度(感度)の問題

(1) 応答表面の傾斜がゆるく(等高線密度が小さい)、シプレックスのサイズが小さい時、最小/極小値に達する速度が遅くなる

これはシプレックスの感度がにぶる事を意味し、効率的かつスピーディな最適化を妨げる要因となる。感度が鈍いのか、鋭いのかという判断を行う事は通常は困難であり、経験的に感知するしかない。しかも、このシプレックスの感度はこの他にも解析に用いるデータセットのデータ分布状態にも大きく影響を受ける。この為、感度が鈍いと判明した時は、鈍さを支配する要因が何に起因するかを把握する事が必要である。

(2) 感度に対する対策: 原因がシプレックス手法に起因する時

① 応答関数を変え、より感度の高いものを利用する。

② シプレックスの大きさを変える

□ シプレックス法等に利用される応答関数について

応答関数はシプレックスの応答面を決定するもので、シプレックスの感度や最適値の発見のし易さ等に影響する極めて重要な関数である。従って、解析が成功するか否かはこの応答関数の適否に大きく左右されると言っても過言ではない。個々の解析作業/目的/扱う数値データの形式等を考慮して最適な関数を採用する事が大事である。

この応答関数の決定にあたり重要な事は、用いる関数が仕事の目的あるいは特性を適切に反映している事である。応答関数に最低必要な条件は、望ましくない所では大きな値を持ち、望ましい所では小さな値を持つような関数を選ぶことである。この反対の時シプレックスは望ましくないものを探し出すことになる。

現在幾つかの関数がこの応答関数として利用されているが、クラス分類問題によく用いられる関数の内主なものを以下に示す。

$$\text{誤差関数: } \sum_i^2 (W - X)$$

$$\text{パーセプトロン関数: } \sum_i |W - X|$$

$$\text{ハイパータンジェント関数: } \sum_i \tan h | \mathbf{W} \cdot \mathbf{X} |$$

なお、式中 i は誤分類されたもののパターンをしめす。従って、ここで示された応答関数は誤分類パターンについてのみ算出される。この値が大きい程誤差が多い事となり、この値が小さくなるようにシプレックスが移動する。

必ずしもここで示された関数に縛られる必要はない。この他にも自分の仕事の目的に応じ、より適切な応答関数を作成する事がシプレックスによる解析に成功するポイントである。

最近では分子力学の歪みエネルギーを最適化する手法として、従来の最小二乗法に変わり得るものとして利用され始めている。

□ シプレックスと応答関数との関係 (実際の問題への適用について)

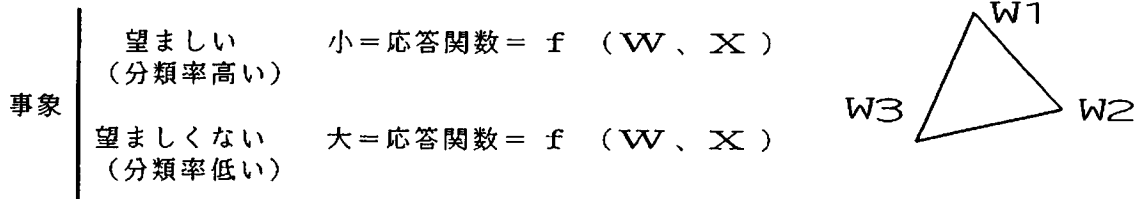
シプレックス法は、最適化対象となる変数群とシプレックスの座標を対応させ、このシプレックスを最適化の程度をモニターする応答関数上を移動させて最適領域を探し出し、この時のシプレックスの位置情報に従った変数値群を最終結果とする最適化手法である。しかし、この手法を自分の最適化問題とリンクする時に大きなギャップを感じる人が多いかもしれない。ここでは、シプレックス法を実際に利用する時に必要となるシプレックスと応答関数との関係について簡単にまとめ、具体例について検討する。

ウェイトベクトル \mathbf{W} による応答関数とシプレックスとの関係

応答関数とシプレックスとの関係はウェイトベクトル \mathbf{W} にて関係つけられる。シプレックスの反転はこのウェイトベクトルの変更という形で実現されている。

一方、自分が解析しようとする問題と応答関数との間はパターンベクトル \mathbf{X} で関係付けられている。この応答関数の値は小さい値である程望ましい結果であり、大きい値である程望ましく無い結果であるというように、自分が求めようとする結果とリンクしている事が必要である。つまり、応答関数はシプレックスを動かす時のモニターとして用いられる。以下に2つの事例を用いて簡単に説明する。

事例1: 分類問題の時 (最適判別関数の選出)



事象と応答関数との関係
分類率に反比例

応答関数とシプレックスとの関係
ウェイトベクトル \mathbf{W}

図6-5にシプレックス法による最適化の流れ図を示す。

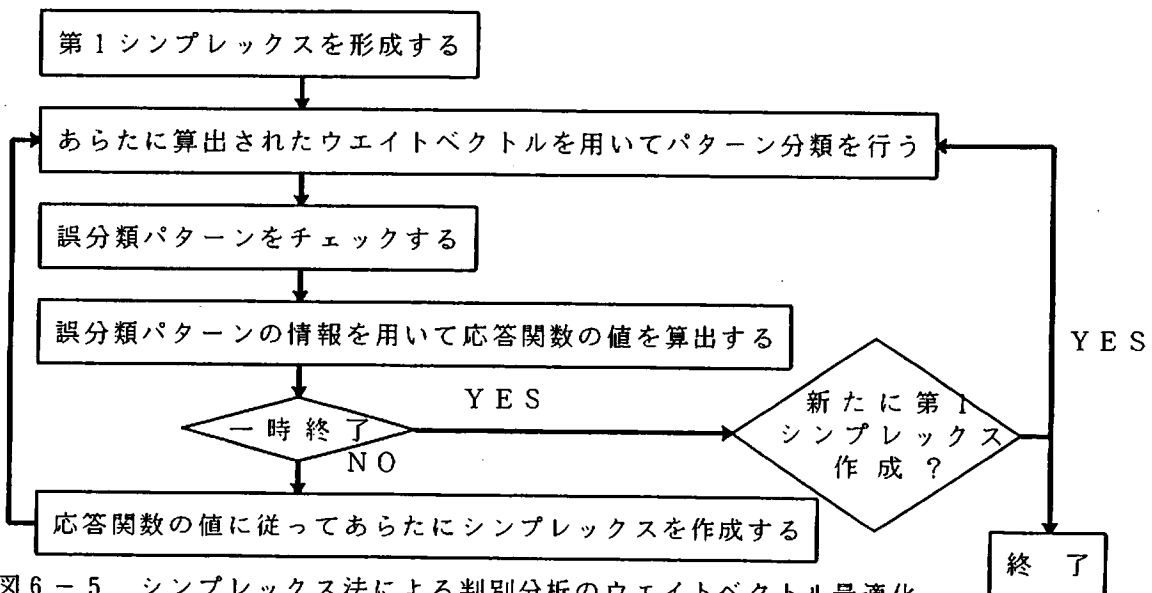


図6-5 シプレックス法による判別分析のウェイトベクトル最適化