

データ解析入門 2 <主成分分析によるデータの可視化>

キーワード：主成分分析、スコアプロット、多変量解析、可視化、次元削減

はじめに

主成分分析 (PCA) は、マーケティングや化学分析などの様々な分野で用いられる多変量解析手法です。一般に、高次元データの特徴をそのまま理解することは困難ですが、PCA では高次元データの特徴的な変動の抽出を図ります。本稿では、PCA の基本的なコンセプトを解説します。

多変量データの取り扱い

ここでは、ワインの化学分析に関するデータセットを PCA により解析します。図 1 にデータセットの一部を示しますが、全 178 サンプルのワインについて、アルコール濃度などの 13 項目が評価されています。このような多変量データは行列で表現でき、178 行×13 列のデータ行列となります。このとき、各変数 ($x_1 \sim x_{13}$) は横方向に並び、13 次元空間における各座標軸に対応します。13 次元空間におけるデータの様子を想像することは難しいため、まずは便宜上選んだ 2 つの変数 x_1 、 x_3 の散布図を確認してみます。図 2(a) に示すように、2 つの座標軸からなる 2 次元平面に 178 個のデータ点がプロットされています。座標軸が 13 個になると、もはや図示できませんが、2 変数の場合と同様に考え、178 個のデータ点が 13 次元空間にプロットされている様子をイメージします。このように考えると、データ行列の各列には $x_1 \sim x_{13}$ の各座標軸上の値 (スコア) が格納されていると理解できます。

分析項目

| Sample ID | x_1 | x_2 | x_3 | x_4 | x_5 | x_6 | x_7 | x_8 | x_9 | x_{10} | x_{11} | x_{12} | x_{13} |
|-----------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|----------|
| 0 | 14.2 | 1.71 | 2.43 | 15.6 | 127 | 2.8 | 3.06 | 0.28 | 2.29 | 5.64 | 1.04 | 3.92 | 1065 |
| 1 | 13.2 | 1.78 | 2.14 | 11.2 | 100 | 2.65 | 2.76 | 0.26 | 1.28 | 4.38 | 1.05 | 3.4 | 1050 |
| 2 | 13.2 | 2.36 | 2.67 | 18.6 | 101 | 2.8 | 3.24 | 0.3 | 2.81 | 5.68 | 1.03 | 3.17 | 1185 |
| 3 | 14.4 | 1.95 | 2.5 | 16.8 | 113 | 3.85 | 3.49 | 0.24 | 2.18 | 7.8 | 0.86 | 3.45 | 1480 |
| 4 | 13.2 | 2.59 | 2.87 | 21 | 118 | 2.8 | 2.69 | 0.39 | 1.82 | 4.32 | 1.04 | 2.93 | 735 |
| ⋮ | | | | | | | | | | | | | |
| 177 | 14.1 | 4.1 | 2.74 | 24.5 | 96 | 2.05 | 0.76 | 0.56 | 1.35 | 9.2 | 0.61 | 1.6 | 560 |

↓

データ行列
(178行×13列)

図 1 データ行列の概念図

データの前処理

解析手法の性質上、PCA ではデータの前処理として中心化を行います。中心化とは、データ行列の各要素から、対応する変数の平均値を引く操作のことであり (式(1))、中心化によりデータ点が原点周りに移動します (図 2(b))。また、今回のデータでは、変数 $x_1 \sim x_{13}$ の単位およびスケールはそれぞれ異なります。そこで、中心化を実行後、データを各変数の標準偏差で割ることで、変数のスケールを統一します (図 2(c))。この操作を標準化といいます (式(2))。標準化によりスケールの違いに由来する影響を解析結果から取り除くことができます。そのため、PCA に限らず、変数の単位が異なる場合などでは、データの標準化がしばしば重要になります。

$$x_j \leftarrow x_j - \bar{x}_j \quad (1)$$

$$x_j \leftarrow \frac{x_j - \bar{x}_j}{\sigma_j} \quad (2)$$

なお、添え字 j は変数 x のインデックス、 \bar{x}_j は変数 x_j の平均値、 σ_j は変数 x_j の標準偏差を表します。

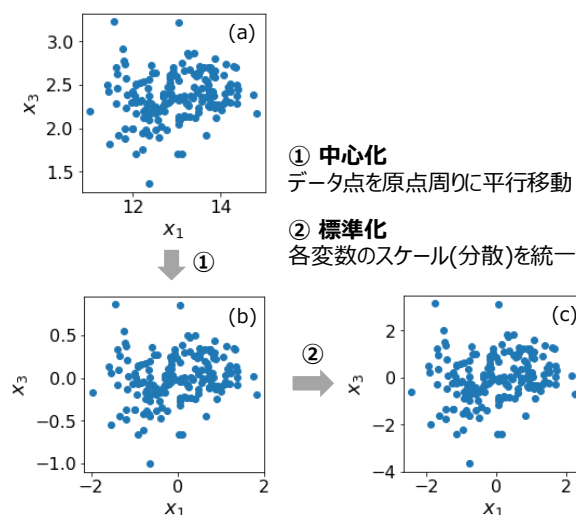


図 2 変数 x_1 および変数 x_3 の散布図：
(a) 前処理なし、(b) 中心化後、(c) 標準化後

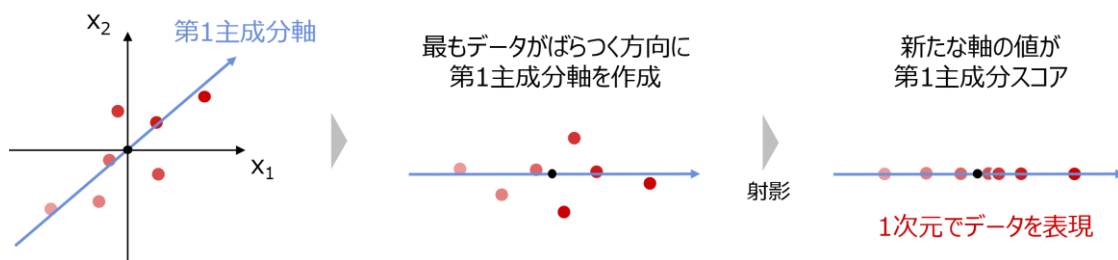


図3 2次元平面におけるPCAによる次元削減

PCAのコンセプト

ここまでで、13次元という多変量データを図示できないことを確認しました。一方で、2変数であれば散布図を作成できました。ただし、図2では全13変数のうち、2変数をプロットしているにすぎず、データ全体の特徴をとらえていません。

PCAでは、変数 $x_1 \sim x_{13}$ の線形結合で表される新たな軸(第1主成分軸)を作成し、多変量データの特徴を少ない変数でとらえられるようにします。つまり、 i 番目のワインの第1主成分スコア $t_{i,1}$ は、

$$t_{i,1} = p_{1,1}x_{i,1} + p_{2,1}x_{i,2} + \dots + p_{13,1}x_{i,13} \quad (3)$$

と表されます。ここで、 $x_{i,j}$ ($j = 1, 2, \dots, 13$) は i 番目のサンプルの変数 x_j の数値であり、 $p_{j,1}$ は後述の処理から求まる重みです。なお、ここでいう主成分(Principal component)とは、PCAで抽出する新たな変数のことであり、化学組成などにおける「主成分」という意味ではありません。

次に、PCAによる主成分軸の求め方を説明します。まず、第1主成分のスコアが最も大きく変動するような第1主成分を見つけだし、その軸への情報の縮約を図ります。簡単のため、2変数から1変数への次元削減を考えます(図3)。今、求めたい第1主成分軸の方向は、第1主成分スコアが最もばらつく方向です。そこで、第1主成分軸方向のデータのばらつきを最大化します。これにより第1主成分軸に元データの主要な変動をできるだけ写しとる操作が実行されます。

2次元平面における次元削減はこれで終了ですが、多変数を扱う場合は、第1主成分軸に直交する中で、次にスコアが最もばらつく第2主成分軸を決定します。続いて、第1および第2主成分軸に直交する中で、スコアが最もばらつく軸を第3主成分軸とする、という操作を繰り返すことで、全ての主成分軸が得られます。

スコアプロットによる可視化

上述の操作により、第1～第13主成分軸が求まります。PCAにより得られる第1および第2主成分スコアの散布図(スコアプロット)を示します(図

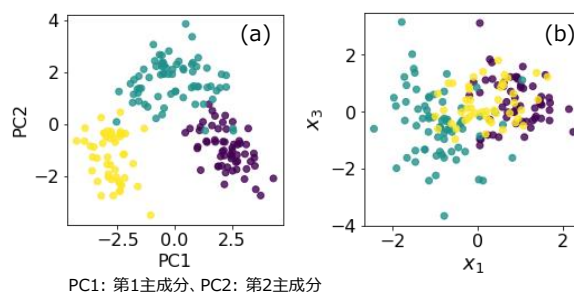


図4 (a) スコアプロット、(b) 図2(c)を色分けした散布図

4(a))。冒頭では述べませんでしたが、本データセットのワインの栽培者はそれぞれ異なり、栽培者ごとに3つのクラスラベルが割り振られています。そこで、図4では3つのクラスごとにプロットを色分けしました。PCAにより得られたスコアプロットでは、クラスごとにデータが集まっていますが(図4(a))、ただ単に変数 x_1, x_3 をプロットするだけでは、データの特徴を十分にとらえられていません(図4(b))。色分けがされていなければ、3種類のワインの存在を見出すことは困難です(図2(c))。このように、PCAによる次元削減を実行することで、データの特徴を可視化できる場合があります。また、可視化に用いる主成分の組み合わせを変更することで、より示唆に富む特徴が見つかることもあります。

おわりに

本稿では、データ解析のための前処理やPCAの概念について紹介しました。PCAは有用な線形次元削減手法であり、データの特徴を確認するためにしばしば用いられます。次稿では、本稿で用いたワインのデータセットを題材に、PCAによるデータマイニングについて紹介します。

参考文献

- 1) UCI Machine Learning Repository, <https://archive.ics.uci.edu/ml/machine-learning-database/s/wine/wine.data> (accessed on April 28th, 2021)