

## 主成分分析

多変量データの持つ構造をより少数個の指標に「圧縮」

○変量の個数を減らすことに伴う、情報の損失はなるべく小さくしたい

○少数変量を利用した分析や可視化(2・3次元の場合)が実現可能

学習データ

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{im}) \in \mathbb{R}^m$$

平均 (ベクトル)

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$$

データ行列

$$\bar{X} = (\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}})^T \in \mathbb{R}^{n \times m}$$

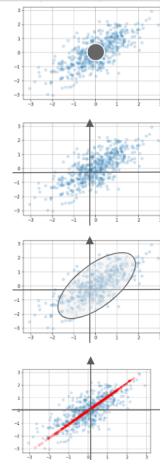
分散共分散行列  
(復習用)

$$\Sigma = \text{Var}(\bar{X}) = \frac{1}{n} \bar{X}^T \bar{X}$$

線形変換後の  
ベクトル

$$\mathbf{s}_j = (s_{1j}, \dots, s_{nj})^T = \bar{X} \mathbf{a}_j \quad \mathbf{a}_j \in \mathbb{R}^m$$

※jは射影軸のインデックス

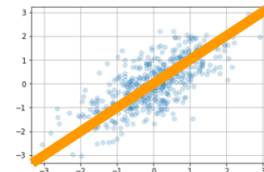
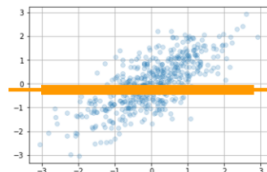
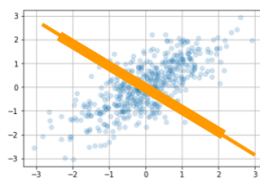


係数ベクトルが変われば線形変換後の値が変化

「情報量の分散の大きさ」ととらえる。

線形変換後の変数の分散が最大となる射影軸を探索

$$\mathbf{s}_j = (s_{1j}, \dots, s_{nj})^T = \bar{X} \mathbf{a}_j \quad \mathbf{a}_j \in \mathbb{R}^m$$



線形変換後の分散

$$\text{Var}(\mathbf{s}_j) = \frac{1}{n} \mathbf{s}_j^T \mathbf{s}_j = \frac{1}{n} (\bar{X} \mathbf{a}_j)^T (\bar{X} \mathbf{a}_j) = \frac{1}{n} \mathbf{a}_j^T \bar{X}^T \bar{X} \mathbf{a}_j = \mathbf{a}_j^T \text{Var}(\bar{X}) \mathbf{a}_j$$

↑↑↑一番右側のグラフが「分散最大」な射影軸がえられる。。

よって、以下の制約付き最適化問題を解けばよい。

制約:  $\mathbf{a}_j^T \mathbf{a}_j = 1 \leftarrow$  ノルムが1.

このような制約を入れないと、例えば、(1, 1) (2, 2) など、大きさが違うベクトルが大量に出てきてしまい解けない。

目的関数

$$\arg \max_{\mathbf{a} \in \mathbb{R}^m} \mathbf{a}_j^T \text{Var}(\bar{X}) \mathbf{a}_j$$

制約条件

$$\mathbf{a}_j^T \mathbf{a}_j = 1$$

ラグランジュ乗数

ラグランジュ関数

$$E(\mathbf{a}_j) = \mathbf{a}_j^T \text{Var}(\bar{X}) \mathbf{a}_j - \lambda (\mathbf{a}_j^T \mathbf{a}_j - 1)$$

目的関数

制約条件

## ラグランジュの未定乗数法

$g(x, y) = 0$ のもとで  $f(x, y)$  を最大化したいという等式制約付きの問題において

$L(x, y, \lambda) = f(x, y) - \lambda g(x, y)$ を作ると

$(\alpha, \beta)$ が極値を与える $\rightarrow (\alpha, \beta)$ は $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} = \frac{\partial L}{\partial \lambda} = 0$ の解。

寄与率: 第1-k主成分において保持している情報量の割合:  $\lambda_k / \sum_{i=1-k} \lambda_i$

大量データの絞り込みなどに利用。

分散共分散行列  
を計算

固有値問題を解く

(最大)m個の固有値と  
固有ベクトルのペアが出現

k番目の固有値(昇順)  
並べ、対応する固有ベクトル  
を第k主成分と呼ぶ

## ● 寄与率

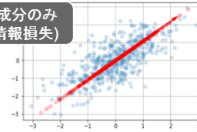
- 第1~元次元分の主成分の分散は、元のデータの分散と一致
  - 2次元のデータを2次元の主成分で表示した時、固有値の和と元のデータの分散が一致
  - 第k主成分の分散は主成分に対応する固有値

$$V_{total} = \sum_{i=1}^m \lambda_i$$

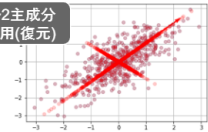
元データの  
総分散

主成分の  
総分散

第1主成分のみ  
利用(情報損失)



第1~2主成分  
を利用(復元)



- 寄与率: 第k主成分の分散の全分散に対する割合 (第k主成分が持つ情報量の割合)
- 累積寄与率: 第1-k主成分まで圧縮した際の情報損失量の割合

$$c_k = \frac{\lambda_k}{\sum_{i=1}^m \lambda_i}$$

第k主成分の  
分散

主成分の  
総分散

$$r_k = \frac{\sum_{j=1}^k \lambda_j}{\sum_{i=1}^m \lambda_i}$$

第1~k主成  
分の分散

主成分の  
総分散