

クレジット:

Mathematics and Informatics Center 文科系のための線形代数・解析Ⅱ
2020 藤堂 眞治・松尾 泰・藤原 毅夫

ライセンス:

利用者は、本講義資料を、教育的な目的に限ってページ単位で利用することができます。特に記載のない限り、本講義資料はページ単位でクリエイティブ・コモンズ 表示-非営利-改変禁止 ライセンスの下に提供されています。

<http://creativecommons.org/licenses/by-nc-nd/4.0/>

本講義資料内には、東京大学が第三者より許諾を得て利用している画像等や、各種ライセンスによって提供されている画像等が含まれています。個々の画像等を本講義資料から切り離して利用することはできません。個々の画像等の利用については、それぞれの権利者の定めるところに従ってください。



§ 主成分分析

多次元 データ $\vec{x}^{(i)} = \begin{bmatrix} x_1^{(i)} \\ \vdots \\ x_m^{(i)} \end{bmatrix} \in \mathbb{R}^m \dots \vec{x}^{(j)} = \begin{bmatrix} x_1^{(j)} \\ \vdots \\ x_m^{(j)} \end{bmatrix} \in \mathbb{R}^m$

(例) 入学試験

$x_i^{(j)}$: 国語の点数. $x_i^{(j)}$: 数学の点数, ... $i=1 \sim n$ 受験生 ID

主成分分析: 多次元データから情報をとける次元を減らしながら 情報の収約 を行う

例) 入学試験: ^{重み}合計点

$$y_i = \sum_{a=1}^p w^{(a)} x_i^{(a)}$$

一つの指標で優劣を測る

重み: 文系: 国語. 理系: 数学

◦ 肥満度: $BMI = (\text{体重}) / (\text{身長})^2$: 一つのデータで判定

$$\log(BMI) = \underbrace{1 \times \log(\text{体重}) - 2 \log(\text{身長})}_{\text{重み}}$$

"重み" をどのようにつけるのか?

統計学より

データの1つの場合

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

平均 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

偏差 $x_i - \bar{x}$

Var(x)

分散 $\sigma = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$
 $= \frac{1}{n} \sum_{i=1}^n x_i^2 - (\bar{x})^2$

複数のデータがある場合

$$x_i^{(a)} \quad a=1, \dots, p, \quad i=1, \dots, n$$

各データの平均

$$\bar{x}^{(a)} = \frac{1}{n} \sum_{i=1}^n x_i^{(a)}$$

分散 $\sigma_{aa} = \frac{1}{n} \sum_{i=1}^n (x_i^{(a)} - \bar{x}^{(a)})^2 = \frac{1}{n} \sum_{i=1}^n (x_i^{(a)})^2 - (\bar{x}^{(a)})^2$
 $\equiv \text{Var}(x^{(a)})$

共分散 $\sigma_{ab} = \frac{1}{n} \sum_{i=1}^n x_i^{(a)} x_i^{(b)} = \text{Cov}(x^{(a)}, x^{(b)})$

相関係数 $\rho_{ab} = \sigma_{ab} / \sqrt{\sigma_{aa} \sigma_{bb}}$

$-1 \leq \rho_{ab} \leq 1$ $\rho_{ab} > 0$ 正の相関 $\rho_{ab} < 0$ 負の相関

$$X = [x_i^{(a)} - \bar{x}^{(a)}]_{\substack{1 \leq i \leq n, \\ 1 \leq a \leq p}} \quad n \times p$$

各データの偏差を
行列の形にまとめたもの

${}^t X \cdot X$: グラム行列.

$$({}^t X \cdot X)_{ab} = \sum_{i=1}^n (x_i^{(a)} - \bar{x}^{(a)}) (x_i^{(b)} - \bar{x}^{(b)}) = \begin{cases} n \text{Var}(x^{(a)}) & : a=b \\ n \text{Cov}(x^{(a)}, x^{(b)}) & : a \neq b \end{cases}$$

グラム行列の固有値方程式

\vec{w} : $p \times 1$ ベクトル

$${}^t X X \cdot \vec{w}^{(a)} = \lambda^{(a)} \vec{w}^{(a)}$$

$w^{(a)}$: 固有ベクトル

$\lambda^{(a)}$: 固有値

* $\lambda^{(a)}$ は 0 以上の実数

$$\begin{pmatrix} w_1^{(a)} \\ \vdots \\ w_n^{(a)} \end{pmatrix}$$

* $\lambda^{(a)} \neq \lambda^{(b)}$ のとき $\vec{w}^{(a)} \perp \vec{w}^{(b)}$

⊙ 固有値方程式の転置は $({}^t AB) = B^t A$ に注意する

$${}^t \vec{w}^{(a)} {}^t X X = \lambda^{(a)} {}^t \vec{w}^{(a)}$$

$${}^t \vec{w}^{(a)} {}^t X X w^{(b)} = {}^t \vec{w}^{(a)} ({}^t X X w^{(b)}) = \lambda^{(b)} {}^t \vec{w}^{(a)} w^{(b)}$$

$$= ({}^t \vec{w}^{(a)} {}^t X X) w^{(b)} = \lambda^{(a)} {}^t \vec{w}^{(a)} w^{(b)}$$

$$(\lambda_a - \lambda_b) {}^t \vec{w}^{(a)} \cdot \vec{w}^{(b)} = 0 \quad \Leftrightarrow \quad \lambda_a \neq \lambda_b \text{ のとき } {}^t \vec{w}^{(a)} \cdot \vec{w}^{(b)} = 0$$

* $\vec{w}^{(a)} = \frac{1}{\sqrt{w^{(a)T} w^{(a)}}} \vec{w}^{(a)}$ とおくと $\vec{w}^{(a)} \cdot \vec{w}^{(b)} = \delta_{ab}$ 正規直交基底

$\vec{w}^{(a)}$ を改めて $\vec{w}^{(a)}$ と書く

* $y_i^{(a)} = \sum_{b=1}^p (x_i^{(b)} - \bar{x}^{(b)}) w_b^{(a)}$ データの重み和

各 a に対して $\vec{y}^{(a)} = \begin{bmatrix} y_1^{(a)} \\ \vdots \\ y_n^{(a)} \end{bmatrix}$ の平均は 0, 分散は $\frac{1}{n} \lambda^{(a)}$

$a \neq b$ ($\lambda^{(a)} \neq \lambda^{(b)}$ を仮定する) $\text{Cov}(y^{(a)}, y^{(b)}) = 0$

☺ $\vec{y}^{(a)} \cdot \vec{y}^{(b)} = \underbrace{w^{(a)T} X \cdot X}_{\lambda_a} w^{(b)} = \lambda_a w^{(a)} \cdot w^{(b)} = \lambda_a \delta_{ab}$

$\text{Var}(y^{(a)}) = \frac{1}{n} \sum_{i=1}^n (y_i^{(a)})^2 = \frac{1}{n} \vec{y}^{(a)} \cdot \vec{y}^{(a)} = \frac{1}{n} \lambda_a$

◎ 固有値を $\lambda^{(1)} > \lambda^{(2)} > \dots > \lambda^{(n)} > 0$ のように順番に並べる

主成分 $\vec{y}^{(1)}$: 最も分散が大きいデータ

第2成分 $\vec{y}^{(2)}$: 分散が2番目に大きいデータ

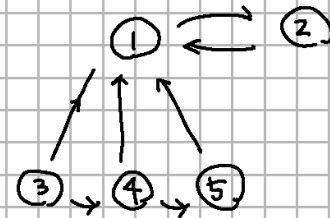
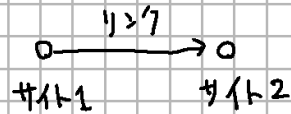
§ Google 行列と Page rank

インターネット 検索 エンジン

- 条件にあてはまるサイトの中で重要なものほど上に出すおに工夫している

"重要性"の目安 : Page rank

※ ネットではユーザは各サイトのリンクをたどってサイトを移動している



リンクは左の方向
だけアットラムで表記
できる

※ どのサイトが最も読まれるのか?

- 多くのサイトからリンクされている
- 重要なサイトからリンクされるとより重要になる

※ リンクを表す行列 $A = [a_{ij}]$ \rightarrow 各目のサイト n_j へのリンク i_1, \dots, i_{n_j}

$$a_{ij} = \begin{cases} \frac{1}{n_j} \\ 0 \end{cases} \quad i = [i_1, \dots, i_{n_j}]$$

上のグラフの場合は

$$A = \begin{bmatrix} 0 & 1 & \frac{1}{2} & \frac{1}{2} & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} & 0 \end{bmatrix}$$

重要性の目安

N-ページ

$$\vec{p} = \begin{bmatrix} p_1 \\ \vdots \\ p_N \end{bmatrix}$$

$N \times 1$ ベクトル

N : サイトの数

$$p_i = \sum_j A_{ij} p_j$$

i 番目のページの価値 = $\sum_{\text{リンク元ページ } N \rightarrow i} \frac{\text{リンク元ページの価値}}{\text{リンク元ページのリンク数}}$

行列表示: $\vec{p} = A\vec{p}$: \vec{p} は A の固有値 1 に対する固有ベクトル

規格化 $\sum_{i=1}^N p_i = 1, p_i \geq 0$

例としてあげた行列の場合

$$\vec{p} = \begin{bmatrix} 1/2 \\ 1/2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

1 と 2 は同等の価値

3, 4, 5 は 価値が 0

極端すぎます?

偶然 ③④⑤ を見るとはありうる..

Google 行列 $G = [g_{ij}]$

$$g_{ij} = \alpha A_{ij} + (1 - \alpha) \frac{1}{N}$$

ランダムにサイトを訪れる効果

例 $\alpha = 0.8$

$$P = \begin{bmatrix} 0.35 \\ 0.32 \\ 0.04 \\ 0.18 \\ 0.11 \end{bmatrix}$$

Google $\alpha \sim 0.85$

Google 行列の性質

- 各行列要素は正
- 任意の j に対し $\sum_{i=1}^N a_{ij} = 1$

固有値 1 の固有ベクトルが存在し unique であること

↑
Perron - Frobenius の定理

▷ 行列要素が全て正の行列 A に対して

▷ $\rho(A) > 0$ の固有値 r

A の固有値 λ に対し $|\lambda| < r$
他

▷ 固有値 r に対する固有ベクトル ($\rho(A)$ ベクトル) の成分は全て正

▷ それ以外の固有値に対する固有ベクトルは負の成分を持つ

Google 行列の場合

• \vec{p} の成分は全て正

• 固有値 1

• \vec{p} 以外の固有ベクトルは負の成分を持つ

\vec{p} は $\rho(A)$ ベクトル

Google 行列のサイズ N : 全世界のサイト数 17兆? 超巨大行列

全ての固有ベクトル, 固有値を求めたい 現実的でない

(固有値を求め方程式 $\det(\lambda I - G) = 0$ 項の数 $N!$)

Perron-Frobenius の定理を用いて少い計算量で \vec{p} のみを得たい.

固有値方程式 $G \vec{u}_i = \lambda_i \vec{u}_i$ $\lambda_1 = 1$ $\vec{u}_1 = \vec{p}$ $|\lambda_i| < 1$ $i \neq 1$

$\Lambda = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_n \end{pmatrix}$ 対角行列 $U = [\vec{u}_1, \dots, \vec{u}_n]$ 固有ベクトルを並べたもの

$$GU = U\Lambda \Rightarrow G = U\Lambda U^{-1}$$

* 任意のベクトル $\vec{v} = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$ $v_i > 0$ $\sum_{i=1}^n v_i = 0$

∴ 対して

$$\lim_{n \rightarrow \infty} G^n \vec{v} = \vec{p}$$

☺

$$G^n = (U \Lambda U^{-1})^n = U \Lambda \cancel{U^{-1}} U \Lambda U^{-1} \cdots \cancel{U^{-1}} U \Lambda U^{-1} = U \Lambda^n U^{-1}$$

$$\rightarrow \Lambda = \begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_N \end{pmatrix} \quad |\lambda_i| < 1 \quad \forall i \text{ の こと}$$

$$\lim_{n \rightarrow \infty} \Lambda^n = \lim_{n \rightarrow \infty} \begin{pmatrix} \lambda_1^n & & \\ & \lambda_2^n & \\ & & \ddots \\ & & & \lambda_N^n \end{pmatrix} = \begin{pmatrix} 1 & & \\ & 0 & \\ & & \ddots \\ & & & 0 \end{pmatrix}$$

$$U = \begin{bmatrix} \vec{u}_1 & \vec{u}_2 & \cdots & \vec{u}_n \\ \parallel \\ \vec{p} \end{bmatrix} \quad U^{-1} = \begin{bmatrix} \vec{w}_1 \\ \vdots \\ \vec{w}_n \end{bmatrix} \quad \text{と } \vec{w}_i \perp \vec{p}$$

$$\lim_{n \rightarrow \infty} G^n = \begin{bmatrix} \vec{u}_1 & \cdots & \vec{u}_n \end{bmatrix} \begin{bmatrix} 1 & & \\ & 0 & \\ & & \ddots \\ & & & 0 \end{bmatrix} \begin{bmatrix} \vec{w}_1 \\ \vdots \\ \vec{w}_n \end{bmatrix} = \begin{matrix} \vec{u}_1 & \vec{w}_1 \\ \uparrow & \uparrow \\ N \times 1 & 1 \times N \text{ 行列} \end{matrix}$$

$$\therefore \lim_{n \rightarrow \infty} G^n \vec{v} = \underbrace{\vec{u}_1}_{= \vec{p}} (\vec{w}_1 \cdot \vec{v}) \propto \vec{p}$$

比例係数が1であること $\sum_{i=1}^N (G \vec{v})_i = \sum_{i,j} G_{ij} v_j = \sum_j v_j \left(\sum_{i=1}^N G_{ij} = 1 \right)$