

主成分分析

基本的な考え方

村田 昇

講義の内容

- 第 1 日 : 主成分分析の考え方
- 第 2 日 : 分析の評価と視覚化

主成分分析の考え方

主成分分析

- 多数の変量のもつ情報の分析・視覚化
 - 変量を効率的に縮約して少数の特徴量を構成する
 - 特徴量に関与する変量間の関係を明らかにする
- PCA (Principal Component Analysis)
 - 構成する特徴量 : **主成分** (principal component)

分析の枠組み

- x_1, \dots, x_p : 変数
- z_1, \dots, z_d : 特徴量 ($d \leq p$)
- 変数と特徴量の関係 (線形結合)

$$z_k = a_{1k}x_1 + \dots + a_{pk}x_p \quad (k = 1, \dots, d)$$

- 特徴量は定数倍の任意性があるので以下を仮定

$$\|a_k\|^2 = \sum_{j=1}^p a_{jk}^2 = 1$$

主成分分析の用語

- 特徴量 z_k
 - 第 k **主成分得点** (principal component score)
 - 第 k **主成分**
- 係数ベクトル a_k
 - 第 k **主成分負荷量** (principal component loading)
 - 第 k **主成分方向** (principal component direction)

分析の目的

- 目的
主成分得点 z_1, \dots, z_d が変数 x_1, \dots, x_p の情報を効率よく反映するように主成分負荷量 a_1, \dots, a_d を観測データから決定する
- 分析の方針 (以下は同値)
 - データの情報を最も保持する変量の **線形結合を構成**
 - データの情報を最も反映する **座標軸を探索**
- **教師なし学習** の代表的手法の 1 つ
 - 特徴抽出: 情報処理に重要な特性を変数に凝集
 - 次元縮約: 入力をできるだけ少ない変数で表現

第 1 主成分の計算

記号の準備

- 変数: x_1, \dots, x_p (p 次元)
- 観測データ: n 個の (x_1, \dots, x_p) の組

$$\{(x_{i1}, \dots, x_{ip})\}_{i=1}^n$$

- ベクトル表現
 - $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^T$: i 番目の観測データ (p 次元空間内の 1 点)
 - $\mathbf{a} = (a_1, \dots, a_p)^T$: 長さ 1 の p 次元ベクトル

係数ベクトルによる射影

- データ \mathbf{x}_i の \mathbf{a} 方向成分の長さ

$$\mathbf{a}^T \mathbf{x}_i \quad (\text{スカラー})$$

- 方向ベクトル \mathbf{a} をもつ直線上への点 \mathbf{x}_i の直交射影

$$(\mathbf{a}^T \mathbf{x}_i) \mathbf{a} \quad (\text{スカラー} \times \text{ベクトル})$$

幾何学的描像

ベクトル \mathbf{a} の選択の指針

- 射影による特徴量の構成
ベクトル \mathbf{a} を **うまく** 選んで観測データ $\mathbf{x}_1, \dots, \mathbf{x}_n$ の情報を最も保持する 1 変量データ z_1, \dots, z_n を構成

$$z_1 = \mathbf{a}^T \mathbf{x}_1, z_2 = \mathbf{a}^T \mathbf{x}_2, \dots, z_n = \mathbf{a}^T \mathbf{x}_n$$

- 特徴量のばらつきの最大化
観測データの **ばらつき** を最も反映するベクトル \mathbf{a} を選択

$$\arg \max_{\mathbf{a}} \sum_{i=1}^n (\mathbf{a}^T \mathbf{x}_i - \mathbf{a}^T \bar{\mathbf{x}})^2, \quad \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i,$$

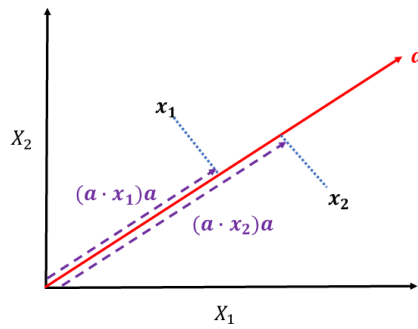


Figure 1: 観測データの直交射影 ($p = 2, n = 2$ の場合)

ベクトル a の最適化

- 最適化問題

制約条件 $\|a\| = 1$ の下で以下の関数を最大化せよ

$$f(a) = \sum_{i=1}^n (a^T x_i - a^T \bar{x})^2$$

- この最大化問題は必ず解をもつ
 - $f(a)$ は連続関数
 - 集合 $\{a \in \mathbb{R}^p : \|a\| = 1\}$ はコンパクト (有界閉集合)

演習

問題

- 以下の問に答えなさい
 - 評価関数 $f(a)$ を以下の中心化したデータ行列で表しなさい

$$X = \begin{pmatrix} x_1^T - \bar{x}^T \\ \vdots \\ x_n^T - \bar{x}^T \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & \cdots & x_{1p} - \bar{x}_p \\ \vdots & & \vdots \\ x_{n1} - \bar{x}_1 & \cdots & x_{np} - \bar{x}_p \end{pmatrix}$$

- 上の結果を用いて次の最適化問題の解の条件を求めなさい

$$\text{maximize } f(a) \quad \text{s.t. } a^T a = 1$$

解答例

- 定義どおりに計算する

$$\begin{aligned} f(a) &= \sum_{i=1}^n (a^T x_i - a^T \bar{x})^2 \\ &= \sum_{i=1}^n (a^T x_i - a^T \bar{x})(x_i^T a - \bar{x}^T a) \\ &= a^T X^T X a \end{aligned}$$

- 回帰分析の Gram 行列を参照
- 制約付き最適化なので未定係数法を用いればよい

$$L(\mathbf{a}, \lambda) = f(\mathbf{a}) + \lambda(1 - \mathbf{a}^T \mathbf{a})$$

の鞍点

$$\frac{\partial}{\partial \mathbf{a}} L(\mathbf{a}, \lambda) = 0$$

を求めればよいので

$$2X^T X \mathbf{a} - 2\lambda \mathbf{a} = 0$$

$$X^T X \mathbf{a} = \lambda \mathbf{a} \quad (\text{固有値問題})$$

第1主成分の解

ベクトル \mathbf{a} の解

- 最適化問題

$$\text{maximize } f(\mathbf{a}) = \mathbf{a}^T X^T X \mathbf{a} \quad \text{s.t. } \mathbf{a}^T \mathbf{a} = 1$$

- 固有値問題

$f(\mathbf{a})$ の極大値を与える \mathbf{a} は $X^T X$ の固有ベクトルとなる

$$X^T X \mathbf{a} = \lambda \mathbf{a}$$

第1主成分

- 固有ベクトル \mathbf{a} に対する $f(\mathbf{a})$ は行列 $X^T X$ の固有値

$$f(\mathbf{a}) = \mathbf{a}^T X^T X \mathbf{a} = \mathbf{a}^T \lambda \mathbf{a} = \lambda$$

- 求める \mathbf{a} は行列 $X^T X$ の最大固有ベクトル (長さ 1)
- **第1主成分負荷量**: 最大 (第一) 固有ベクトル \mathbf{a}
- **第1主成分得点**

$$z_{i1} = a_1 x_{i1} + \cdots + a_p x_{ip} = \mathbf{a}^T \mathbf{x}_i, \quad (i = 1, \dots, n)$$

Gram 行列の性質

Gram 行列の固有値

- $X^T X$ は非負定値対称行列
- $X^T X$ の固有値は 0 以上の実数
 - 固有値を重複を許して降順に並べる

$$\lambda_1 \geq \cdots \geq \lambda_p \quad (\geq 0)$$

- 固有値 λ_k に対する固有ベクトルを \mathbf{a}_k (長さ 1) とする

$$\|\mathbf{a}_k\| = 1, \quad (k = 1, \dots, p)$$

Gram 行列のスペクトル分解

- $\mathbf{a}_1, \dots, \mathbf{a}_p$ は **互いに直交** するようとりとることができる

$$j \neq k \quad \Rightarrow \quad \mathbf{a}_j^\top \mathbf{a}_k = 0$$

- 行列 $X^\top X$ (非負定値対称行列) のスペクトル分解

$$\begin{aligned} X^\top X &= \lambda_1 \mathbf{a}_1 \mathbf{a}_1^\top + \lambda_2 \mathbf{a}_2 \mathbf{a}_2^\top + \cdots + \lambda_p \mathbf{a}_p \mathbf{a}_p^\top \\ &= \sum_{k=1}^p \lambda_k \mathbf{a}_k \mathbf{a}_k^\top \end{aligned}$$

- 固有値と固有ベクトルによる行列の表現

演習

問題

- 以下の問に答えなさい
 - Gram 行列のスペクトル分解において λ_j と \mathbf{a}_j が固有値・固有ベクトルとなることを確かめなさい

$$X^\top X = \sum_{k=1}^p \lambda_k \mathbf{a}_k \mathbf{a}_k^\top$$

- 以下の行列を用いて Gram 行列のスペクトル分解を書き直しなさい

$$A = \begin{pmatrix} \mathbf{a}_1^\top \\ \vdots \\ \mathbf{a}_p^\top \end{pmatrix}, \quad \Lambda = \begin{pmatrix} \lambda_1 & 0 & \cdots & 0 \\ 0 & \lambda_2 & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & \lambda_p \end{pmatrix}$$

解答例

- 固有ベクトルの直交性に注意する

$$\begin{aligned} X^\top X \mathbf{a}_j &= \sum_{k=1}^p \lambda_k \mathbf{a}_k \mathbf{a}_k^\top \mathbf{a}_j && \text{(直交性)} \\ &= \lambda_j \mathbf{a}_j \mathbf{a}_j^\top \mathbf{a}_j && \text{(単位ベクトル)} \\ &= \lambda_j \mathbf{a}_j \end{aligned}$$

- 転置に注意して計算する

$$X^\top X = A^\top \Lambda A$$

第2主成分以降の計算

第2主成分の考え方

- 第1主成分
 - 主成分負荷量: ベクトル \mathbf{a}_1
 - 主成分得点: $\mathbf{a}_1^\top \mathbf{x}_i$ ($i = 1, \dots, n$)
- 第1主成分負荷量に関してデータが有する情報

$$(\mathbf{a}_1^\top \mathbf{x}_i) \mathbf{a}_1 \quad (i = 1, \dots, n)$$

- 第1主成分を取り除いた観測データ (分析対象)

$$\tilde{\mathbf{x}}_i = \mathbf{x}_i - (\mathbf{a}_1^\top \mathbf{x}_i) \mathbf{a}_1 \quad (i = 1, \dots, n)$$

第2主成分の最適化

- 最適化問題
制約条件 $\|\mathbf{a}\| = 1$ の下で以下の関数を最大化せよ

$$\tilde{f}(\mathbf{a}) = \sum_{i=1}^n (\mathbf{a}^\top \tilde{\mathbf{x}}_i - \mathbf{a}^\top \bar{\tilde{\mathbf{x}}})^2 \quad \text{ただし} \quad \bar{\tilde{\mathbf{x}}} = \frac{1}{n} \sum_{i=1}^n \tilde{\mathbf{x}}_i$$

演習

問題

- 以下の問に答えなさい
 - 以下の中心化したデータ行列を X と \mathbf{a}_1 で表しなさい

$$\tilde{X} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top - \bar{\tilde{\mathbf{x}}}^\top \\ \vdots \\ \tilde{\mathbf{x}}_n^\top - \bar{\tilde{\mathbf{x}}}^\top \end{pmatrix}$$

- 上の結果を用いて次の最適化問題の解を求めなさい

$$\text{maximize} \quad \tilde{f}(\mathbf{a}) \quad \text{s.t.} \quad \mathbf{a}^\top \mathbf{a} = 1$$

解答例

- 定義どおりに計算する

$$\tilde{X} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top - \bar{\tilde{\mathbf{x}}}^\top \\ \vdots \\ \tilde{\mathbf{x}}_n^\top - \bar{\tilde{\mathbf{x}}}^\top \end{pmatrix} = X - X \mathbf{a}_1 \mathbf{a}_1^\top$$

- Gram 行列 $\tilde{X}^\top \tilde{X}$ を計算する

$$\begin{aligned}
\tilde{X}^T \tilde{X} &= (X - X\mathbf{a}_1\mathbf{a}_1^T)^T (X - X\mathbf{a}_1\mathbf{a}_1^T) \\
&= X^T X - X^T X\mathbf{a}_1\mathbf{a}_1^T - \mathbf{a}_1\mathbf{a}_1^T X^T X + \mathbf{a}_1\mathbf{a}_1^T X^T X\mathbf{a}_1\mathbf{a}_1^T \\
&= X^T X - \lambda_1 \mathbf{a}_1\mathbf{a}_1^T \\
&= \sum_{k=2}^P \lambda_k \mathbf{a}_k \mathbf{a}_k^T
\end{aligned}$$

元の Gram 行列 $X^T X$ の固有ベクトル \mathbf{a}_1 の固有値が 0 となっていると考えることができる

第 2 主成分以降の解

第 2 主成分

- Gram 行列 $\tilde{X}^T \tilde{X}$ の固有ベクトル \mathbf{a}_1 の固有値は 0

$$\tilde{X}^T \tilde{X} \mathbf{a}_1 = 0$$

- Gram 行列 $\tilde{X}^T \tilde{X}$ の最大固有値は λ_2
- 解は第 2 固有値 λ_2 に対応する固有ベクトル \mathbf{a}_2

-
- 以下同様に第 k 主成分負荷量は $X^T X$ の第 k 固有値 λ_k に対応する固有ベクトル \mathbf{a}_k

解析の事例

データセットについて

- 総務省統計局より取得した都道府県別の社会生活統計指標 (自然環境・経済基盤) の一部
 - 総務省 <https://www.e-stat.go.jp/SG1/estat/List.do?bid=000001083999&cycode=0>
 - データ https://noboru-murata.github.io/multivariate-analysis/data/japan_social.csv
- * Pref: 都道府県名
- * Forest: 森林面積割合 (%) 2014 年
- * Agri: 就業者 1 人当たり農業産出額 (販売農家) (万円) 2014 年
- * Ratio: 全国総人口に占める人口割合 (%) 2015 年
- * Land: 土地生産性 (耕地面積 1 ヘクタール当たり) (万円) 2014 年
- * Goods: 商業年間商品販売額 [卸売業 + 小売業] (事業所当たり) (百万円) 2013 年
- * Area: 地方区分

社会生活統計指標の分析

- データ (の一部) の内容
- データの散布図
- データの箱ひげ図
- 正規化したデータ (の一部)

Table 1: 社会生活統計指標

Pref	Forest	Agri	Ratio	Land	Goods	Area
Hokkaido	67.9	1150.6	4.23	96.8	283.3	Hokkaido
Aomori	63.8	444.7	1.03	186	183	Tohoku
Iwate	74.9	334.3	1.01	155.2	179.4	Tohoku
Miyagi	55.9	299.9	1.84	125.3	365.9	Tohoku
Akita	70.5	268.7	0.81	98.5	153.3	Tohoku
Yamagata	68.7	396.3	0.88	174.1	157.5	Tohoku
Fukushima	67.9	236.4	1.51	127.1	184.5	Tohoku
Ibaraki	31	479	2.3	249.1	204.9	Kanto
Tochigi	53.2	402.6	1.55	199.6	204.3	Kanto
Gumma	63.8	530.6	1.55	321.6	270	Kanto
Saitama	31.9	324.7	5.72	247	244.7	Kanto
Chiba	30.4	565.5	4.9	326.1	219.7	Kanto
Tokyo	34.8	268.5	10.63	404.7	1062.6	Kanto
Kanagawa	38.8	322.8	7.18	396.4	246.1	Kanto
Niigata	63.5	308.6	1.81	141.9	205.5	Chubu
Toyama	56.6	276.1	0.84	98.5	192.4	Chubu
Ishikawa	66	271.3	0.91	112	222.9	Chubu
Fukui	73.9	216.1	0.62	98.5	167.3	Chubu
Yamanashi	77.8	287.4	0.66	325.3	156.2	Chubu
Nagano	75.5	280	1.65	211.3	194.4	Chubu
Gifu	79	283.7	1.6	192.1	167.9	Chubu
Shizuoka	63.1	375.8	2.91	314.5	211.4	Chubu
Aichi	42.2	472.3	5.89	388.9	446.9	Chubu
Mie	64.3	310.6	1.43	174.3	170.1	Kansai
Shiga	50.5	222.8	1.11	104.9	170.7	Kansai
Kyoto	74.2	267.8	2.05	212.5	196.7	Kansai
Osaka	30.1	216.3	6.96	238.8	451.2	Kansai
Hyogo	66.7	261.2	4.35	197.7	212.5	Kansai
Nara	76.8	207	1.07	182.7	147	Kansai
Wakayama	76.4	251.1	0.76	278.4	136.4	Kansai
Tottori	73.3	249.9	0.45	187.6	162.2	Chugoku
Shimane	77.5	214.1	0.55	140.8	141.1	Chugoku
Okayama	68	254.8	1.51	184.9	207.8	Chugoku
Hiroshima	71.8	286.2	2.24	192.2	304.6	Chugoku
Yamaguchi	71.6	216.9	1.11	125.8	158.9	Chugoku
Tokushima	75.2	315.4	0.59	313.5	134.5	Shikoku
Kagawa	46.4	249.5	0.77	242.9	232.9	Shikoku
Ehime	70.3	288.5	1.09	231.6	179.4	Shikoku
Kochi	83.3	354.2	0.57	339.9	137.9	Shikoku
Fukuoka	44.5	381	4.01	255.6	295.7	Kyushu
Saga	45.2	468.7	0.66	230.3	137.9	Kyushu
Nagasaki	58.4	428.9	1.08	296	154	Kyushu
Kumamoto	60.4	456.6	1.41	285.5	172.5	Kyushu
Oita	70.7	360.1	0.92	222.8	148.3	Kyushu
Miyazaki	75.8	739.1	0.87	487.7	170.6	Kyushu
Kagoshima	63.4	736.5	1.3	351.2	169.4	Kyushu
Okinawa	46.1	452.4	1.13	232.8	145.4	Kyushu

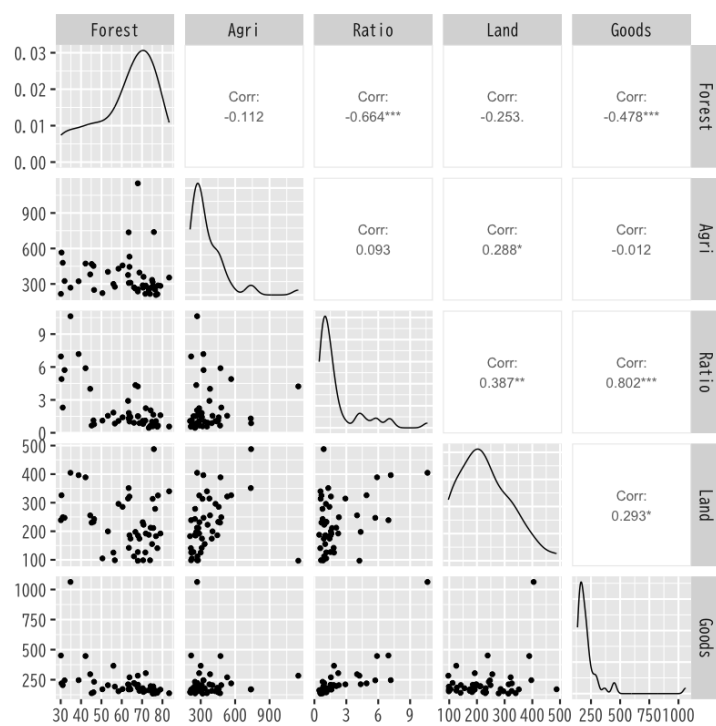


Figure 2: 散布図

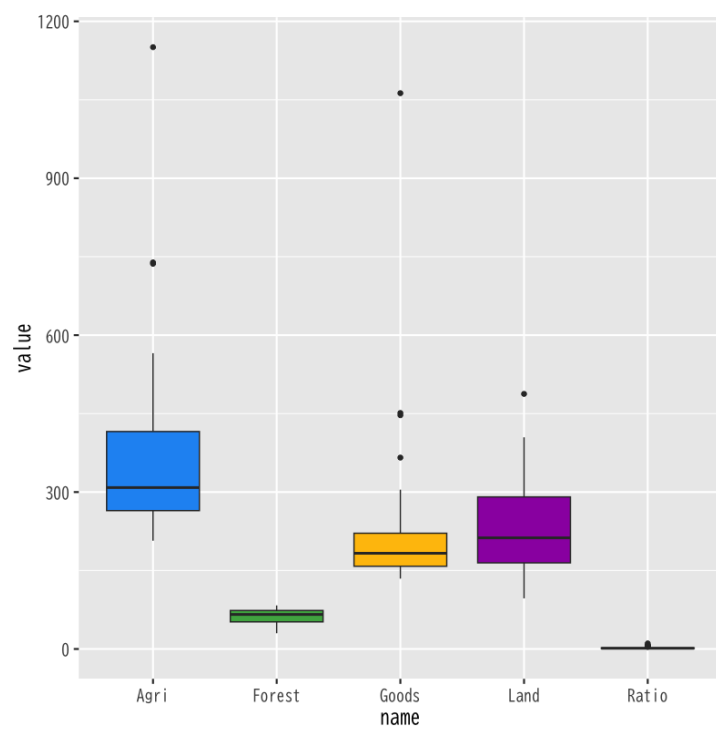


Figure 3: 箱ひげ図

Table 2: 社会生活統計指標

Pref	Forest	Agri	Ratio	Land	Goods	Area
Hokkaido	0.425	4.63	0.979	-1.4	0.421	Hokkaido
Aomori	0.151	0.489	-0.512	-0.446	-0.274	Tohoku
Iwate	0.892	-0.159	-0.521	-0.776	-0.299	Tohoku
Miyagi	-0.376	-0.361	-0.134	-1.1	0.993	Tohoku
Akita	0.599	-0.544	-0.614	-1.38	-0.48	Tohoku
Yamagata	0.479	0.205	-0.581	-0.574	-0.451	Tohoku
Fukushima	0.425	-0.734	-0.288	-1.08	-0.264	Tohoku
Ibaraki	-2.04	0.691	0.0801	0.229	-0.123	Kanto
Tochigi	-0.556	0.242	-0.269	-0.301	-0.127	Kanto
Gumma	0.151	0.994	-0.269	1.01	0.329	Kanto
Saitama	-1.98	-0.215	1.67	0.207	0.153	Kanto
Chiba	-2.08	1.2	1.29	1.05	-0.02	Kanto
Tokyo	-1.78	-0.546	3.96	1.9	5.82	Kanto
Kanagawa	-1.52	-0.227	2.35	1.81	0.163	Kanto
Niigata	0.131	-0.31	-0.148	-0.918	-0.118	Chubu
Toyama	-0.329	-0.501	-0.6	-1.38	-0.209	Chubu
Ishikawa	0.298	-0.529	-0.567	-1.24	0.00214	Chubu
Fukui	0.826	-0.853	-0.703	-1.38	-0.383	Chubu
Yamanashi	1.09	-0.435	-0.684	1.05	-0.46	Chubu
Nagano	0.933	-0.478	-0.223	-0.175	-0.195	Chubu
Gifu	1.17	-0.456	-0.246	-0.381	-0.379	Chubu
Shizuoka	0.105	0.0846	0.364	0.93	-0.0776	Chubu
Aichi	-1.29	0.651	1.75	1.73	1.56	Chubu
Mie	0.185	-0.298	-0.325	-0.572	-0.364	Kansai
Shiga	-0.737	-0.814	-0.474	-1.31	-0.36	Kansai
Kyoto	0.846	-0.55	-0.0364	-0.163	-0.179	Kansai
Osaka	-2.1	-0.852	2.25	0.119	1.58	Kansai
Hyogo	0.345	-0.588	1.04	-0.321	-0.07	Kansai
Nara	1.02	-0.907	-0.493	-0.482	-0.524	Kansai
Wakayama	0.993	-0.648	-0.637	0.543	-0.598	Kansai
Tottori	0.786	-0.655	-0.782	-0.429	-0.419	Chugoku
Shimane	1.07	-0.865	-0.735	-0.93	-0.565	Chugoku
Okayama	0.432	-0.626	-0.288	-0.458	-0.103	Chugoku
Hiroshima	0.686	-0.442	0.0521	-0.38	0.569	Chugoku
Yamaguchi	0.672	-0.849	-0.474	-1.09	-0.442	Chugoku
Tokushima	0.913	-0.27	-0.717	0.919	-0.611	Shikoku
Kagawa	-1.01	-0.657	-0.633	0.163	0.0715	Shikoku
Ehime	0.585	-0.428	-0.484	0.042	-0.299	Shikoku
Kochi	1.45	-0.0422	-0.726	1.2	-0.587	Shikoku
Fukuoka	-1.14	0.115	0.877	0.299	0.507	Kyushu
Saga	-1.09	0.63	-0.684	0.0281	-0.587	Kyushu
Nagasaki	-0.209	0.396	-0.488	0.732	-0.476	Kyushu
Kumamoto	-0.0756	0.559	-0.335	0.619	-0.347	Kyushu
Oita	0.612	-0.0076	-0.563	-0.0522	-0.515	Kyushu
Miyazaki	0.953	2.22	-0.586	2.78	-0.36	Kyushu
Kagoshima	0.125	2.2	-0.386	1.32	-0.369	Kyushu
Okinawa	-1.03	0.534	-0.465	0.0548	-0.535	Kyushu

Table 3

	PC1	PC2	PC3	PC4	PC5
Forest	-0.487	0.105	-0.457	0.686	-0.268
Agri	0.134	0.812	0.479	0.305	0.035
Ratio	0.585	-0.151	0.045	0.164	-0.778
Land	0.355	0.485	-0.742	-0.290	0.069
Goods	0.526	-0.269	-0.095	0.571	0.562

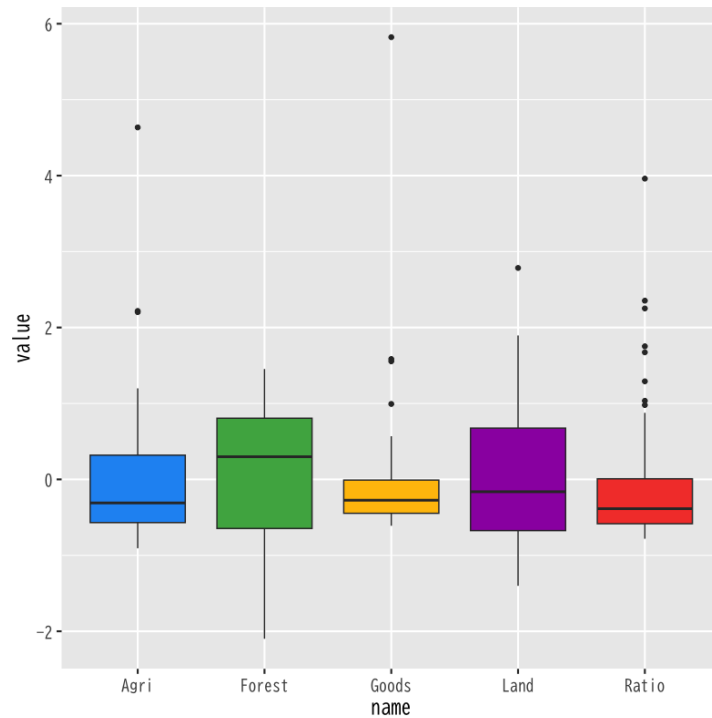


Figure 4: 箱ひげ図 (データを正規化)

- 正規化したデータの箱ひげ図
- 主成分負荷量を計算 (正規化後)
- 主成分方向から読み取れること
 - 第1: 人の多さに関する成分 (正の向きほど人が多い)
 - 第2: 農業生産力に関する成分 (正の向きほど高い)
- 主成分得点の表示

次回の予定

- 第1日: 主成分分析の考え方
- 第2日: 分析の評価と視覚化

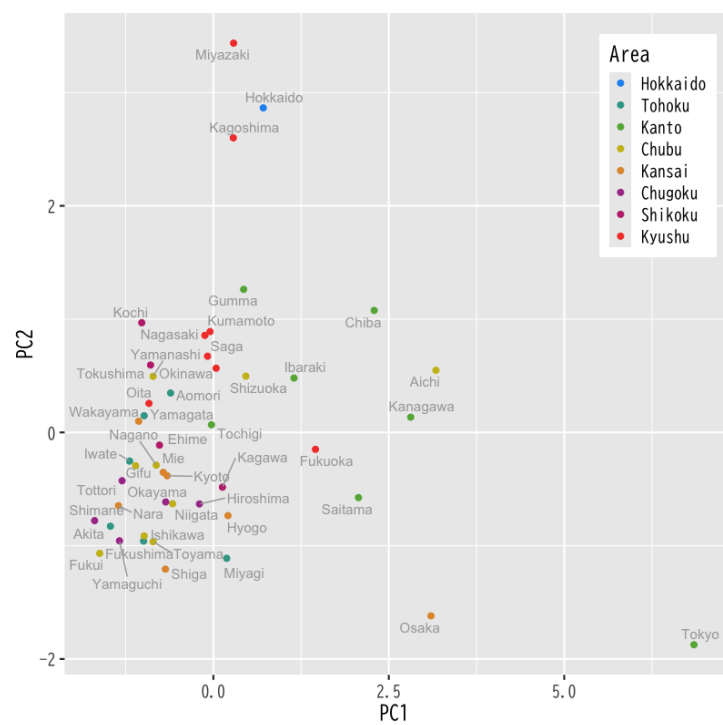


Figure 5: 主成分得点による散布図