

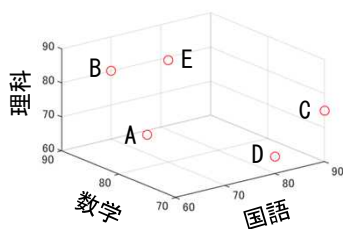
特徴を集約する 主成分分析

Principal Component Analysis

1. PCAで特長を集約
2. 活用例(ネットで検索)
3. PCAの手順(ミスユニバース)
4. PCAの手順(ポケモンGo)
5. PCAと最小二乗法

主成分分析とは？

3次元のデータ



学生	国語	数学	理科
A	70	80	90
B	70	90	80
C	90	70	75
D	80	70	65
E	60	75	75

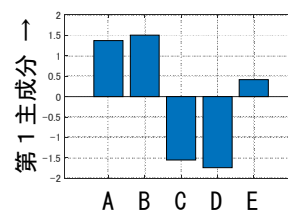
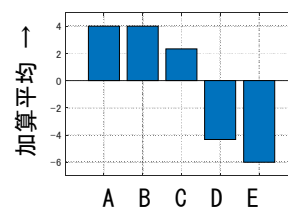
加算平均



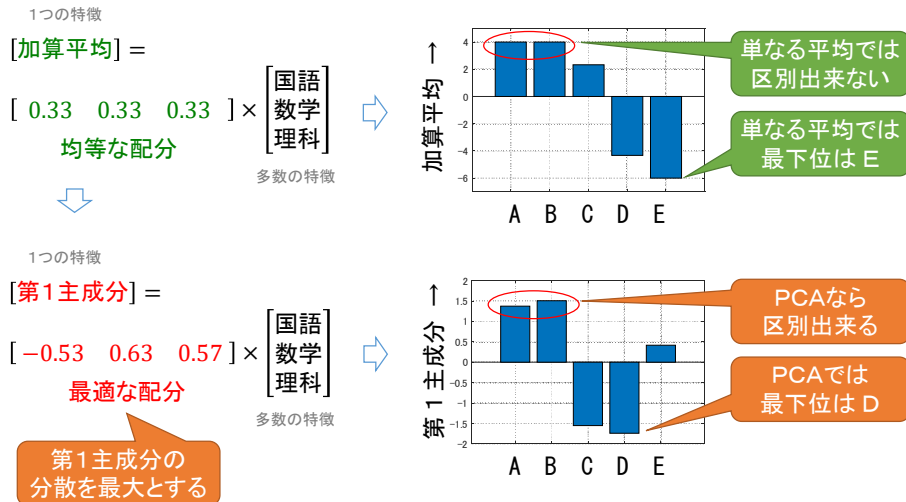
PCA



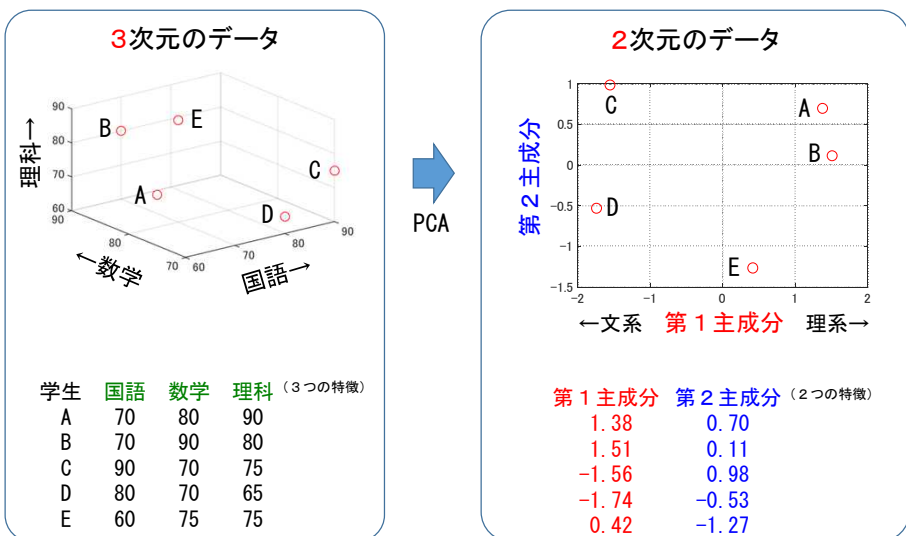
1次元のデータ



最適な配分で 多数の特徴を1つに集約する



多数の特徴を2つに集約



PCAの手順（概要）

学生	国語	数学	理科	← 3つの属性 (多数の特徴)
A	70	80	90	
B	70	90	80	
C	90	70	75	
D	80	70	65	
E	60	75	75	

寄与率

第1主成分 → 65.3 %
第2主成分 → 22.3 %

解釈

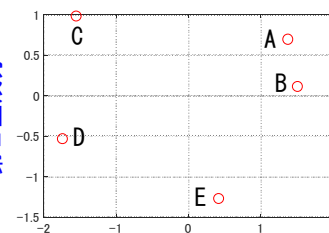
第1主成分 → 理系の力
第2主成分 → 広範な力

PCA

$$\begin{bmatrix} \text{第1主成分} \\ \text{第2主成分} \end{bmatrix} = \begin{bmatrix} -0.53 & 0.63 & 0.57 \\ 0.79 & 0.12 & 0.60 \end{bmatrix} \begin{bmatrix} \text{国語} \\ \text{数学} \\ \text{理科} \end{bmatrix}$$

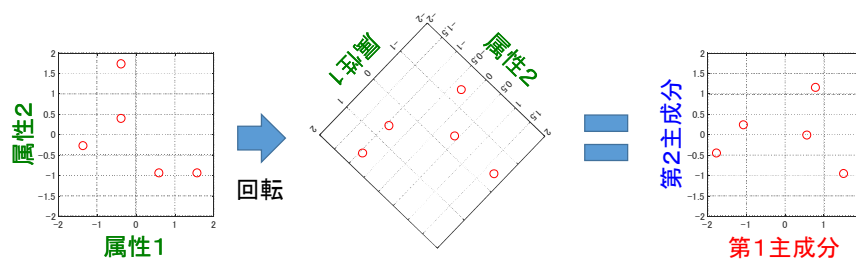
学生	第1主成分	第2主成分
A	1.38	0.70
B	1.51	0.11
C	-1.56	0.98
D	-1.74	-0.53
E	0.42	-1.27

第2主成分



←文系 第1主成分 理系→

PCAは特徴を回転する



$$\begin{bmatrix} \text{第1主成分} \\ \text{第2主成分} \end{bmatrix} = \begin{bmatrix} -0.71 & 0.71 \\ -0.71 & -0.71 \end{bmatrix} \begin{bmatrix} \text{属性1} \\ \text{属性2} \end{bmatrix}$$

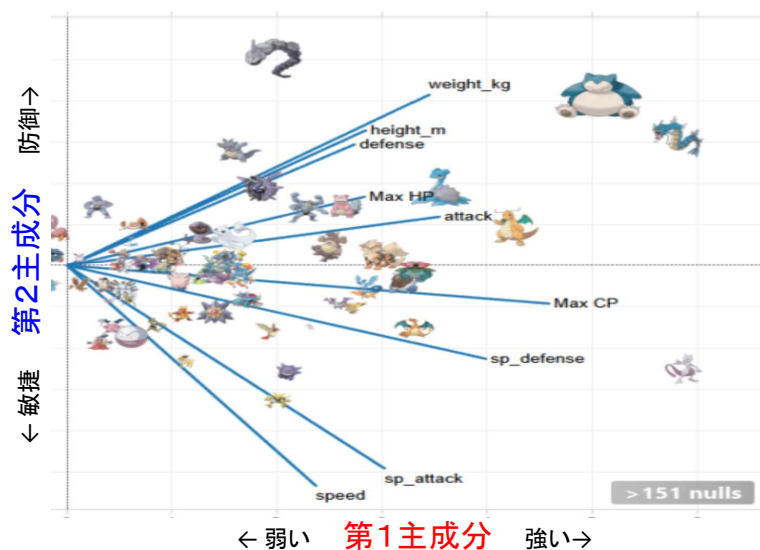
$$= \begin{bmatrix} \cos \frac{-3\pi}{4} & -\sin \frac{-3\pi}{4} \\ \sin \frac{-3\pi}{4} & \cos \frac{-3\pi}{4} \end{bmatrix} \begin{bmatrix} \text{属性1} \\ \text{属性2} \end{bmatrix}$$

特徴を集約する 主成分分析

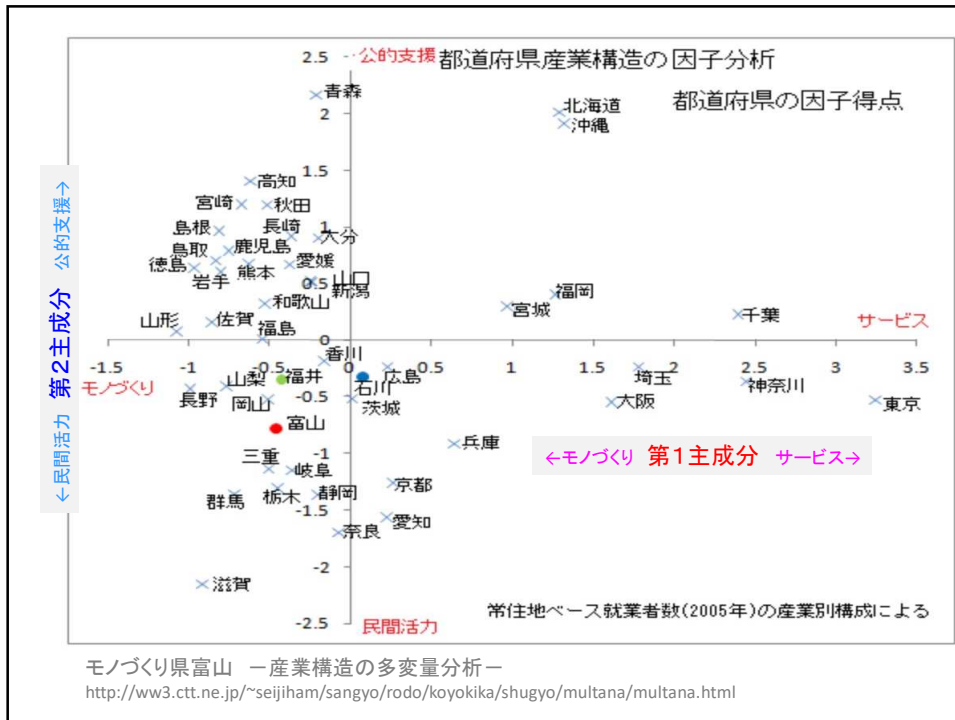
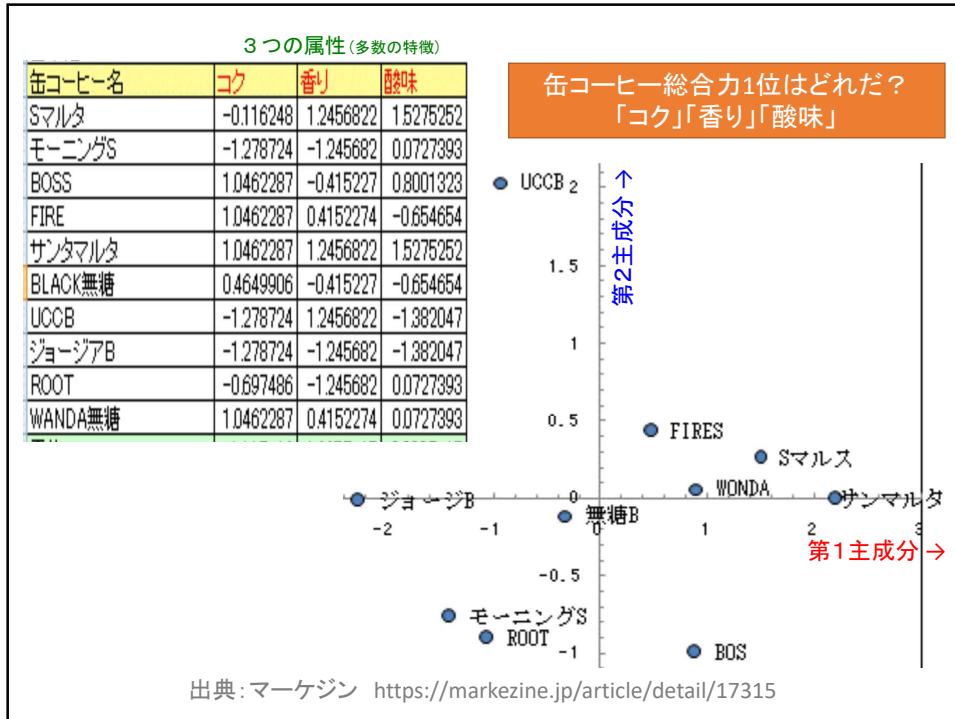
Principal Component Analysis

1. PCAで特長を集約
2. 活用例(ネットで検索)
3. PCAの手順(ミスユニバース)
4. PCAの手順(ポケモンGo)
5. PCAと最小二乗法

種々の属性(戦闘力、体力、体重、身長、...)を主成分分析



出典: 秀和システム, 岩橋智宏ほか, Tableauから始めるデータサイエンス



特徴を集約する 主成分分析

Principal Component Analysis

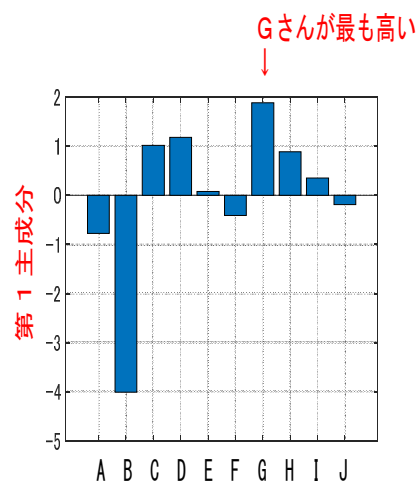
1. PCAで特長を集約
2. 活用例(ネットで検索)
3. PCAの手順(ミスユニバース)
4. PCAの手順(ポケモンGo)
5. PCAと最小二乗法

PCAに出来ること

	身長	体重	V	W	H
Aさん	165	53	86	56	92
Bさん	160	47	84	52	92
Cさん	166	55	86	64	89
Dさん	164	56	90	60	95
Eさん	168	55	87	56	87
Fさん	164	54	87	57	92
Gさん	168	54	94	58	97
Hさん	169	55	88	57	92
Iさん	169	53	86	58	93
Jさん	166	56	84	57	90

5種類の特徴を

多数の属性(多数の特徴)



1つに集約できる

第1主成分

与えられたデータ

	身長	体重	V	W	H	
Aさん	165	53	86	56	92	$= x_1^T$
Bさん	160	47	84	52	92	
Cさん	166	55	86	64	89	
Dさん	164	56	90	60	95	
Eさん	168	55	87	56	87	
Fさん	164	54	87	57	92	
Gさん	168	54	94	58	97	
Hさん	169	55	88	57	92	
Iさん	169	53	86	58	93	
Jさん	166	56	84	57	90	
平均	165.9	53.8	87.2	57.5	91.9	$= \bar{x}^T$
標準偏差	2.66	2.48	2.82	2.91	2.70	$= \sigma^T$

生のデータ
 $= X$

Tは転置

データを標準化する

	身長	体重	V	W	H	
Aさん	-0.34	-0.32	-0.43	-0.52	0.04	
Bさん	-2.22	-2.74	-1.13	-1.89	0.04	
Cさん	0.04	0.48	-0.43	2.24	-1.07	
Dさん	-0.71	0.89	0.99	0.86	1.15	
Eさん	0.79	0.48	-0.07	-0.52	-1.81	
Fさん	-0.71	0.08	-0.07	-0.17	0.04	
Gさん	0.79	0.08	2.41	0.17	1.89	
Hさん	1.16	0.48	0.28	-0.17	0.04	
Iさん	1.16	-0.32	-0.43	0.17	0.41	
Jさん	0.04	0.89	-1.13	-0.17	-0.70	$= \tilde{x}_n^T$
平均	165.9	53.8	87.2	57.5	91.9	$= \bar{x}^T$
標準偏差	2.66	2.48	2.82	2.91	2.70	$= \sigma^T$

標準化されたデータ
 $= \tilde{X}$

$$\tilde{x}_i = \frac{x_i - \bar{x}}{\sigma}$$

データの標準化

←標準化する前の平均

←標準化する前の標準偏差

分散共分散行列を計算

	身長	体重	V	W	H
身長 →	1.00	0.62	0.36	0.37	-0.04
体重 →	0.61	1.00	0.33	0.67	-0.11
V →	0.36	0.33	1.00	0.28	0.69
W →	0.37	0.67	0.28	1.00	-0.01
H →	-0.04	-0.11	0.69	-0.01	1.00

$$= \frac{1}{n} \tilde{X}^T \tilde{X} = C_{xx}$$

分散共分散行列



C_{xx} に対する固有値問題

$$C_{xx}U = UL$$



問題を解く

固有値と固有ベクトルを計算

$$C_{xx}U = UL$$



[U, L] = eig(Cxx);
MATLABのコマンド

$L =$
固有値

0.18	0.00	0.00	0.00	0.00
0.00	0.26	0.00	0.00	0.00
0.00	0.00	0.64	0.00	0.00
0.00	0.00	0.00	1.54	0.00
0.00	0.00	0.00	0.00	2.37

$U =$
固有
ベクトル

第2主成分の係数ベクトル

0.10	0.46	-0.71	-0.18	0.49
0.37	-0.69	0.06	-0.29	0.55
-0.68	-0.24	-0.09	0.53	0.44
-0.11	0.48	0.69	-0.21	0.49
0.62	0.15	0.11	0.75	0.16

第1主成分の係数ベクトル

第1主成分を計算

	身長 ↓	体重 ↓	V ↓	W ↓	H ↓		第1主成分の係数ベクトル ↓		
Aさん→	-0.34	-0.32	-0.43	-0.52	0.04	×	<div style="border: 1px solid red; padding: 2px;"> 0.49 0.55 0.44 0.49 0.16 </div>	=	<div style="border: 1px solid red; padding: 2px;"> -0.78 -4.01 1.01 1.18 0.08 -0.41 1.88 0.88 0.35 -0.19 </div>
Bさん→	-2.22	-2.74	-1.13	-1.89	0.04				
Cさん→	0.04	0.48	-0.43	2.24	-1.07				
Dさん→	-0.71	0.89	0.99	0.86	1.15				
Eさん→	0.79	0.48	-0.07	-0.52	-1.81				
Fさん→	-0.71	0.08	-0.07	-0.17	0.04				
Gさん→	0.79	0.08	2.41	0.17	1.89				
Hさん→	1.16	0.48	0.28	-0.17	0.04				
Iさん→	1.16	-0.32	-0.43	0.17	0.41				
Jさん→	0.04	0.89	-1.13	-0.17	-0.70				

= \bar{X}
標準化されたデータ

5種類の特長が

→

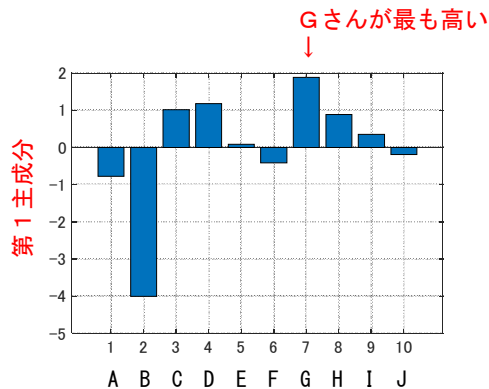
↑
第1主成分
1つに集約された

第1主成分を表示

AからJはミスユニバース日本代表

Aさん→	-0.78
Bさん→	-4.01
Cさん→	1.01
Dさん→	1.18
Eさん→	0.08
Fさん→	-0.41
Gさん→	1.88
Hさん→	0.88
Iさん→	0.35
Jさん→	-0.19

↑
第1主成分

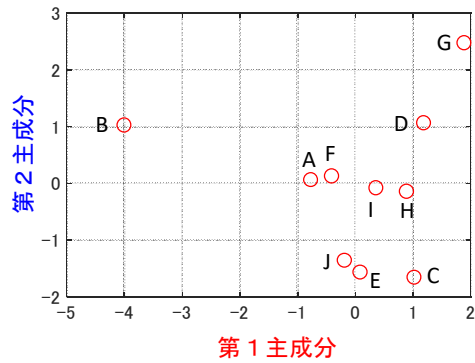


Gさん
昭和34年ミスユニバース
世界大会で優勝！

第2主成分も計算

Aさん→	-0.78	0.07
Bさん→	-4.01	1.03
Cさん→	1.01	-1.65
Dさん→	1.18	1.07
Eさん→	0.08	-1.56
Fさん→	-0.41	0.13
Gさん→	1.88	2.48
Hさん→	0.88	-0.14
Iさん→	0.35	-0.08
Jさん→	-0.19	-1.35

↑ ↑
第1主成分 第2主成分



鈴木義一郎, 情報量基準による統計解析入門, 講談社サイエンティフィック

各成分の意味

第1主成分の係数ベクトル
↓

	身長	体重	V	W	H	
Aさん→	-0.34	-0.32	-0.43	-0.52	0.04	×
	0.49	0.55	0.44	0.49	0.16	
	←身長	←体重	←V	←W	←H	
						高く 重く 大きく 大きく ---

全体的な体格の良さ ←

第2主成分の係数ベクトル
↓

	身長	体重	V	W	H	
Aさん→	-0.34	-0.32	-0.43	-0.52	0.04	×
	-0.18	-0.29	0.53	-0.21	0.75	
	←身長	←体重	←V	←W	←H	
						--- 軽く 大きく 小さく 大きく

プロポーションの良さ ←

PCAで特徴量を集約できた 各成分の寄与率

AからJは、
ミスユニバース日本代表

データは、
体型に関する5種類の属性

↓ PCA

第1主成分の1種類だけでも
47.5%を説明できる

第1と2主成分の2種類なら
78.3%を説明できる

第1~3主成分の3種類なら
91.1%を説明できる

寄与率 (%)

3.6 5.2 12.8 30.8 47.5



$L =$
固有値

0.18	0.00	0.00	0.00	0.00
0.00	0.26	0.00	0.00	0.00
0.00	0.00	0.64	0.00	0.00
0.00	0.00	0.00	1.54	0.00
0.00	0.00	0.00	0.00	2.37

第2主成分の係数ベクトル

$U =$
固有
ベクトル

0.10	0.46	-0.71	-0.18	0.49
0.37	-0.69	0.06	-0.29	0.55
-0.68	-0.24	-0.09	0.53	0.44
-0.11	0.48	0.69	-0.21	0.49
0.62	0.15	0.11	0.75	0.16

↑
第1主成分の係数ベクトル

特徴を集約する 主成分分析

Principal Component Analysis

1. PCAで特長を集約
2. 活用例(ネットで検索)
3. PCAの手順(ミスユニバース)
4. PCAの手順(ポケモンGo)
5. PCAと最小二乗法

MATLAB

MATLAB

① 元のデータ

```

close all; clear all;
%-----元のデータ
X=[
    165    53    86    56    92;
    160    47    84    52    92;
    166    55    86    64    89;
    164    56    90    60    95;
    168    55    87    56    87;
    164    54    87    57    92;
    168    54    94    58    97;
    169    55    88    57    92;
    169    53    86    58    93;
    166    56    84    57    90
];

```

MATLAB

② PCAの計算

```

[n1, n2]=size(X);
%-----データを標準化
%する
xm=mean(X);
% xs=std(X, 0); % N-1 で正規化
xs=std(X, 1); % N で正規化
Xt=(X-xm)./xs;

%-----分散共分散行列
%を計算
Cxx=Xt'*Xt;
Cxx=Cxx/n1;

%-----固有値問題
[U, L] = eig(Cxx);

%-----主成分
Y=Xt*U;

```

MATLAB

③ 計算の結果

```

%----- 各成分の寄与率
fprintf("寄与率\n");
for j=1:n2
    fprintf("%5.1f ", L(j, j)/sum(L(:))*100);
end; fprintf("\n");

%----- 図示：第1主成分と第2主成分
figure('Position', [010 010 300 200]);
plot(Y(:, n2), Y(:, n2-1), 'or'); grid on;

%----- 図示：元のデータ
figure('Position', [010 210 300 200]);
plot3(X(:, 1), X(:, 2), X(:, 3), 'or'); grid on;

%----- 図示：第1主成分
figure('Position', [310 010 300 200]);
bar(Y(:, n2)); grid on;

```

MATLAB

①~③全体

```

235 193 310;
239 261 193;
330 190 169;
189 280 215;
264 226 190;
487 60 128
];
end
[n1, n2]=size(X);
%----- データを標準化
する
xm=mean(X);
% xs=std(X, 0); % N-1 で正規化
xs=std(X, 1); % N で正規化
Xt=(X-xm)./xs;
%----- 分散共分散行列
を計算
Cxx=Xt'*Xt;
Cxx=Cxx/n1;
%----- 固有値問題
[U, L] = eig(Cxx);
%----- 主成分
Y=Xt*U;
%----- 各成分の寄与率
fprintf("寄与率\n");

```

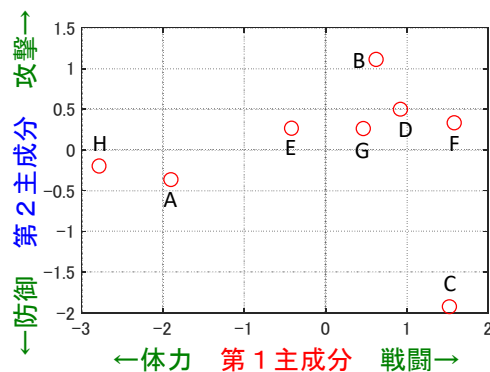
特徴を集約する 主成分分析

Principal Component Analysis

1. PCAで特長を集約
2. 活用例(ネットで検索)
3. PCAの手順(ミスユニバース)
4. PCAの手順(ポケモンGo)
5. PCAと最小二乗法

PCAに出来ること

		体力	攻撃	防御
A	ハピナス	496	129	169
B	ケッキング	284	290	166
C	ルギア	235	193	310
D	ガブリアス	239	261	193
E	カビゴン	330	190	169
F	パルキア	189	280	215
G	メルメタル	264	226	190
H	ラッキー	487	60	128



3種類の特長を



2種類に集約できる

与えられたデータ

	体力	攻撃	防御	
A. ハピナス	496	129	169	$= x_1^T$
B. ケッキング	284	290	166	
C. ルギア	235	193	310	
D. ガブリアス	239	261	193	
E. カビゴン	330	190	169	
F. パルキア	189	280	215	
G. メルメタル	264	226	190	
H. ラッキー	487	60	128	

生のデータ
 $= X$

平均	315.5	203.6	192.5	$= \bar{x}^T$
標準偏差	108.5	73.8	50.4	$= \sigma^T$

データを標準化する

	体力	攻撃	防御	
A. ハピナス	1.66	-1.01	-0.47	
B. ケッキング	-0.29	1.17	-0.53	
C. ルギア	-0.74	-0.14	2.33	
D. ガブリアス	-0.71	0.78	0.01	
E. カビゴン	0.13	-0.18	-0.47	
F. パルキア	-1.17	1.03	0.45	
G. メルメタル	-0.47	0.30	-0.05	
H. ラッキー	1.58	-1.95	-1.28	$= \tilde{x}_n^T$

標準化されたデータ
 $= \tilde{X}$

標準化前の平均	315.5	203.6	192.5	$= \bar{x}^T$
標準化前の標準偏差	108.5	73.8	50.4	$= \sigma^T$

$$\tilde{x}_i = \frac{x_i - \bar{x}}{\sigma}$$

データの標準化

分散共分散行列を計算

	体力	攻撃	防御	
	↓	↓	↓	
体力 →	1.00	-0.86	-0.62	= $\frac{1}{n} \tilde{X}^T \tilde{X} = C_{xx}$
攻撃 →	-0.86	1.00	0.32	
防御 →	-0.62	0.32	1.00	



C_{xx} に対する固有値問題

$$C_{xx}U = UL$$



問題を解く

固有値と固有ベクトルを計算

$$C_{xx}U = UL$$



[U, L] = eig(Cxx);
MATLABのコマンド

$L =$
固有値

0.07	0.00	0.00
0.00	0.70	0.00
0.00	0.00	2.29

$U =$
固有
ベクトル

0.75	-0.11	-0.65
0.59	0.55	0.58
0.29	-0.83	0.48

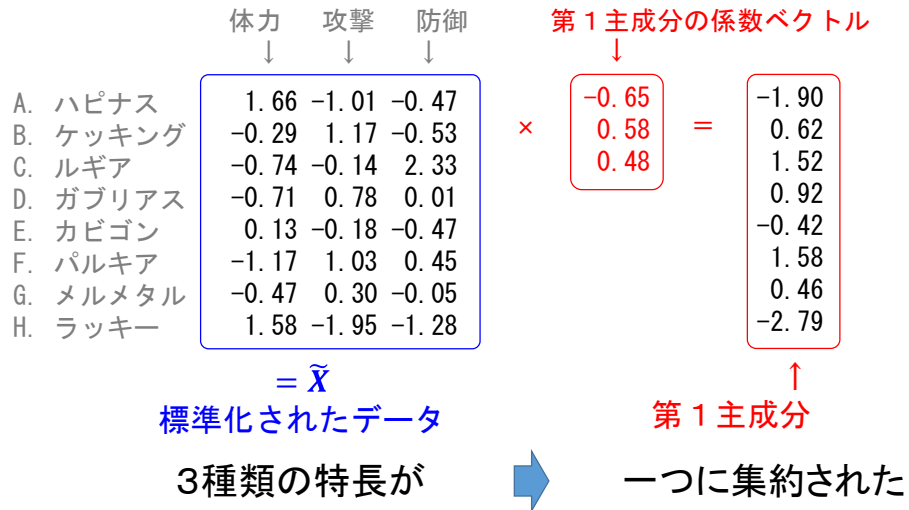
第2主成分の係数ベクトル



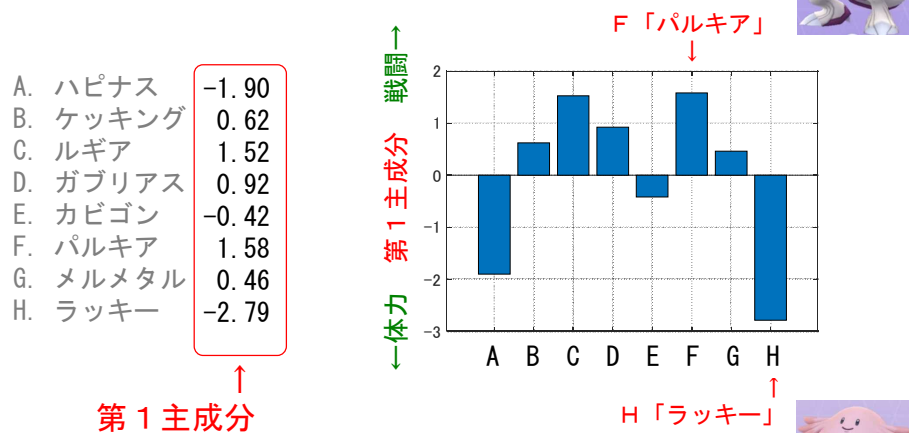
第1主成分の係数ベクトル



第1主成分を計算



第1主成分を表示

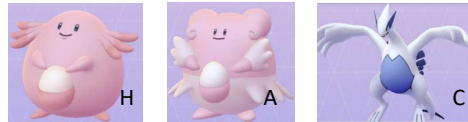
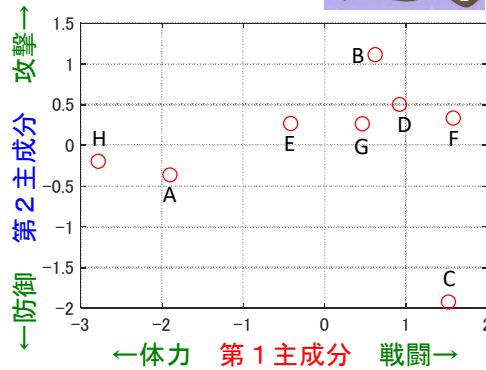


第2主成分も計算



A. ハピナス	-1.90	-0.36
B. ケッキング	0.62	1.11
C. ルギア	1.52	-1.92
D. ガブリアス	0.92	0.50
E. カビゴン	-0.42	0.27
F. パルキア	1.58	0.34
G. メルメタル	0.46	0.26
H. ラッキー	-2.79	-0.20

↑ 第1主成分 ↑ 第2主成分



各成分の意味

第1主成分の係数ベクトル

体力	攻撃	防御	
↓	↓	↓	
ハピナス → 1.66	-1.01	-0.47	×

-0.65 ← 体力 小さく
 0.58 ← 攻撃 大きく
 0.48 ← 防御 大きく

戦闘(攻防)か体力か ←

第2主成分の係数ベクトル

体力	攻撃	防御	
↓	↓	↓	
ハピナス → 1.66	-1.01	-0.47	×

-0.11 ← 体力 ---
 0.55 ← 攻撃 大きく
 -0.83 ← 防御 小さく

攻撃か防御か ←

PCAで特徴量を集約できた 各成分の寄与率

AからHは、
ポケモンのキャラクター

データは、
それぞれ、3種類の属性

↓ PCA

第1主成分の1種類だけでも
74.3%を説明できる

第1と2主成分の2種類なら
97.7%を説明できる

第1~3主成分の3種類なら
100%を説明できる

寄与率 (%)
2.3 23.4 74.3

$L =$
固有値

0.07	0.00	0.00
0.00	0.70	0.00
0.00	0.00	2.29

$U =$
固有
ベクトル

0.75	-0.11	-0.65
0.59	0.55	0.58
0.29	-0.83	0.48

第2主成分の係数ベクトル
↓

↑
第1主成分の係数ベクトル

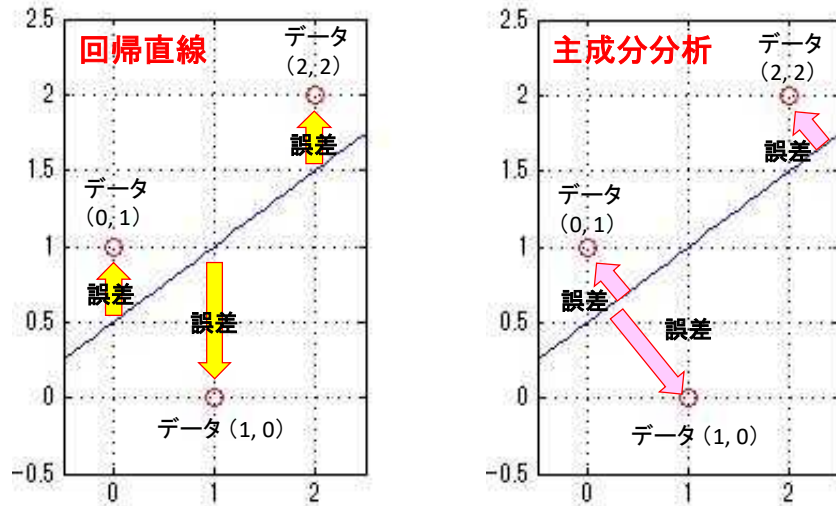
特徴を集約する 主成分分析

Principal Component Analysis

1. PCAで特長を集約
2. 活用例(ネットで検索)
3. PCAの手順(ミスユニバース)
4. PCAの手順(ポケモンGo)
5. PCAと最小二乗法

参考

線形回帰と主成分分析



参考

共分散行列の固有値問題

$$\hat{\mathbf{u}} = \arg \min_{\mathbf{u}} \mathbf{u}^T \mathbf{C}_{xx} \mathbf{u} \quad s. t. \quad \mathbf{u}^T \mathbf{u} = 1$$

$$J(\mathbf{u}) = \mathbf{u}^T \mathbf{C}_{xx} \mathbf{u} - \lambda(\mathbf{u}^T \mathbf{u} - 1)$$

$$\frac{\partial J(\mathbf{u})}{\partial \mathbf{u}} = 2\mathbf{C}_{xx} \mathbf{u} - 2\lambda \mathbf{u} = 0$$

$$\mathbf{C}_{xx} \mathbf{u} = \lambda \mathbf{u}$$

$$\mathbf{C}_{xx} \mathbf{U} = \mathbf{U} \mathbf{L}$$

$$\begin{cases} \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_m) \\ \mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_m) \end{cases}$$

原田達也「画像認識」講談社

参考

主成分の分散を最大化(1/2)

入力データ (n個, 2次元) $\tilde{\mathbf{x}}_i = \begin{bmatrix} \tilde{x}_{i,1} \\ \tilde{x}_{i,2} \end{bmatrix}$ ($i = 1, 2, \dots, n$)



係数ベクトル $\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$ 但し, $u_1^2 + u_2^2 = 1$ とする

主成分 $y_i = [u_1 \quad u_2] \begin{bmatrix} \tilde{x}_{i,1} \\ \tilde{x}_{i,2} \end{bmatrix} = \mathbf{u} \tilde{\mathbf{x}}_i$



主成分の分散 $V = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$ を最大化するように
係数ベクトル \mathbf{u} を決めよ

参考

主成分の分散を最大化(2/2)

主成分の分散 $V = \frac{1}{n} \sum_{i=1}^n (u_1 \tilde{x}_{i,1} + u_2 \tilde{x}_{i,2})^2$ ← 入力データは標準化されているので平均はゼロ

$$= u_1^2 \frac{1}{n} \sum_{i=1}^n \tilde{x}_{i,1}^2 + 2u_1 u_2 \frac{1}{n} \sum_{i=1}^n \tilde{x}_{i,1} \tilde{x}_{i,2} + u_2^2 \frac{1}{n} \sum_{i=1}^n \tilde{x}_{i,2}^2$$

$$= u_1^2 s_{11} + 2u_1 u_2 s_{12} + u_2^2 s_{22}$$



最小二乗法 $(\hat{u}_1, \hat{u}_2) = \arg \min_{u_1, u_2} V \quad \text{s.t.} \quad u_1^2 + u_2^2 = 1$

$$L = u_1^2 s_{11} + 2u_1 u_2 s_{12} + u_2^2 s_{22} - \lambda(u_1^2 + u_2^2 - 1)$$



$$\frac{\partial L}{\partial u_1} = 2u_1 s_{11} + 2u_2 s_{12} - 2\lambda u_1 = 0$$

$$\frac{\partial L}{\partial u_2} = 2u_1 s_{12} + 2u_2 s_{22} - 2\lambda u_2 = 0$$

$$\begin{bmatrix} s_{11} - \lambda & s_{12} \\ s_{21} & s_{22} - \lambda \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix} = \mathbf{C}_{xx} \quad \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} = \mathbf{u}$$



$$(\mathbf{C}_{xx} - \lambda \mathbf{I}) \mathbf{u} = \mathbf{0}$$



固有値問題