

主成分分析



内容

- 主成分分析
 - 主成分分析について
 - 成績データの解析
 - 「R」で主成分分析
 - 相関行列による主成分分析
 - 寄与率・累積寄与率
 - 因子負荷量
 - 主成分得点

主成分分析



次元の縮小と主成分分析

主成分分析

- 次元の縮小に関する手法

□ 次元の縮小

- 国語、数学、理科、社会、英語の総合点
⇒ 5次元データから1次元データへの縮約
- 体形評価：BMI (Body Mass Index) 判定
肥満度の判定方法の1つで、次の式で得られる。

$$\text{BMI} = \frac{\text{体重 (kg)}}{\text{身長 (m)}^2} \Rightarrow \text{2次元データを1次元データに縮約}$$

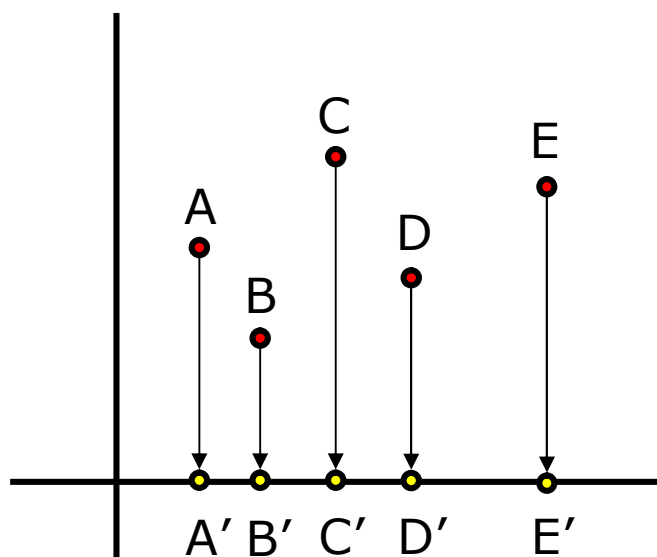
主成分分析とは

□ 主成分分析

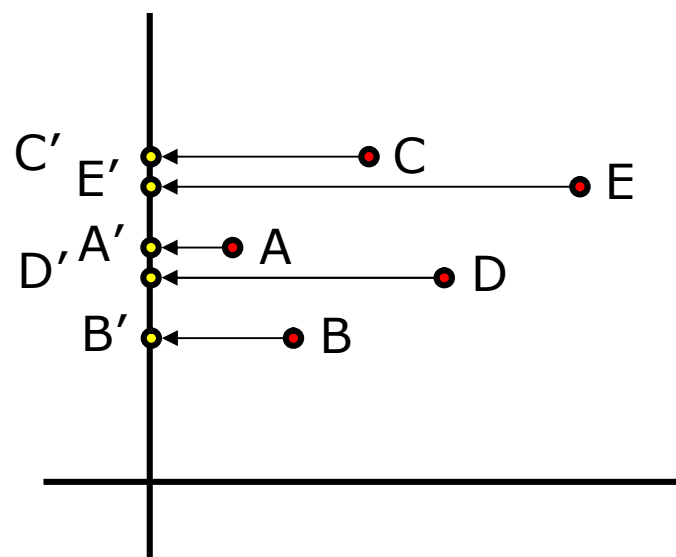
- 多次元データのもつ情報をできるだけ損わずに低次元空間に情報を縮約する方法
- 多次元データを2次元・3次元データに縮約できれば、データ全体の雰囲気を見覚化することができる。視覚化により、データが持つ情報を解釈しやすくなる。

次元の縮約と情報の損失

- 2次元のデータを1次元に縮約することを考える。



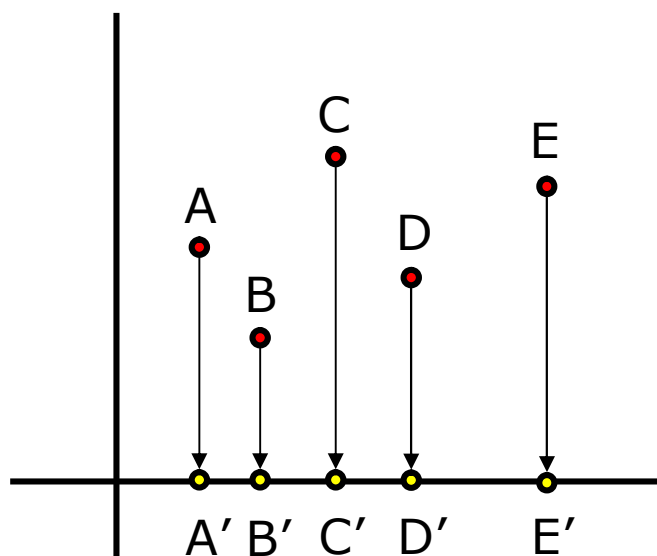
縮約の方法 ①
縦軸の情報の損失



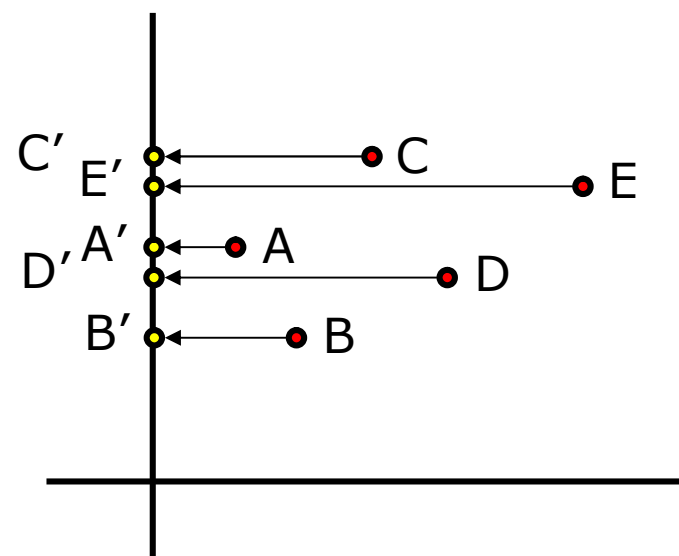
縮約の方法 ②
横軸の情報の損失

情報量と分散

- 射影したデータのバラツキが大きいほど、もとのデータの情報を多く含んでいると考えられる。



個体差が現れやすい

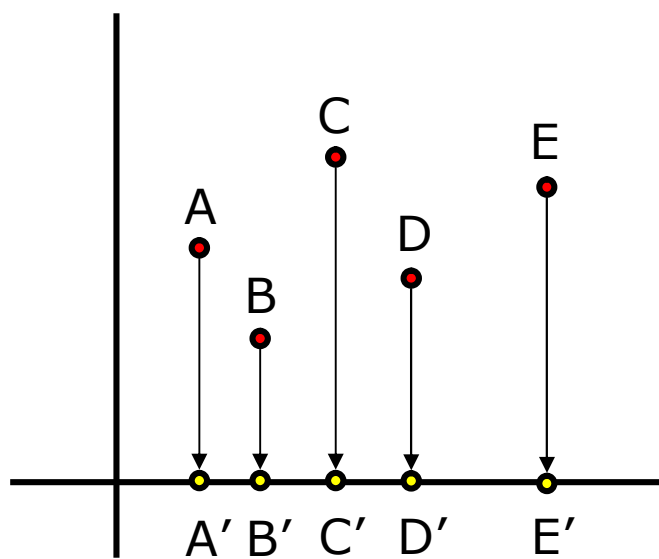


個体差が現れにくい

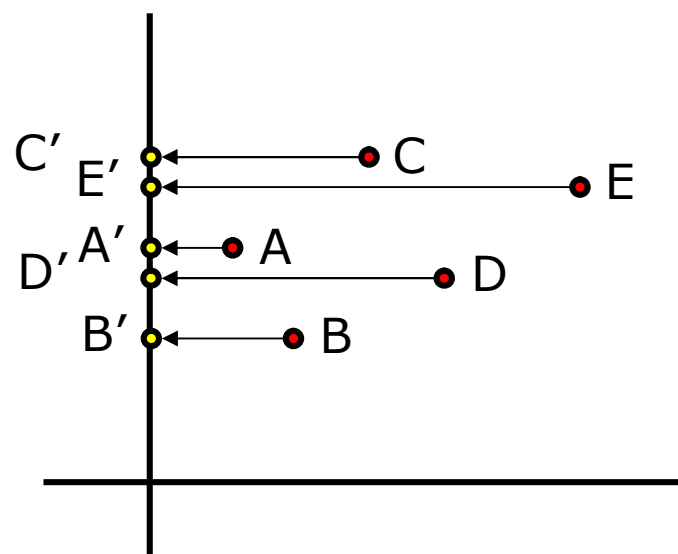
情報量 ↔ 分散

主成分分析の目的

- もとのデータの情報の損失ができるだけ小さくなるような軸を探したい。



情報の損失が少ない

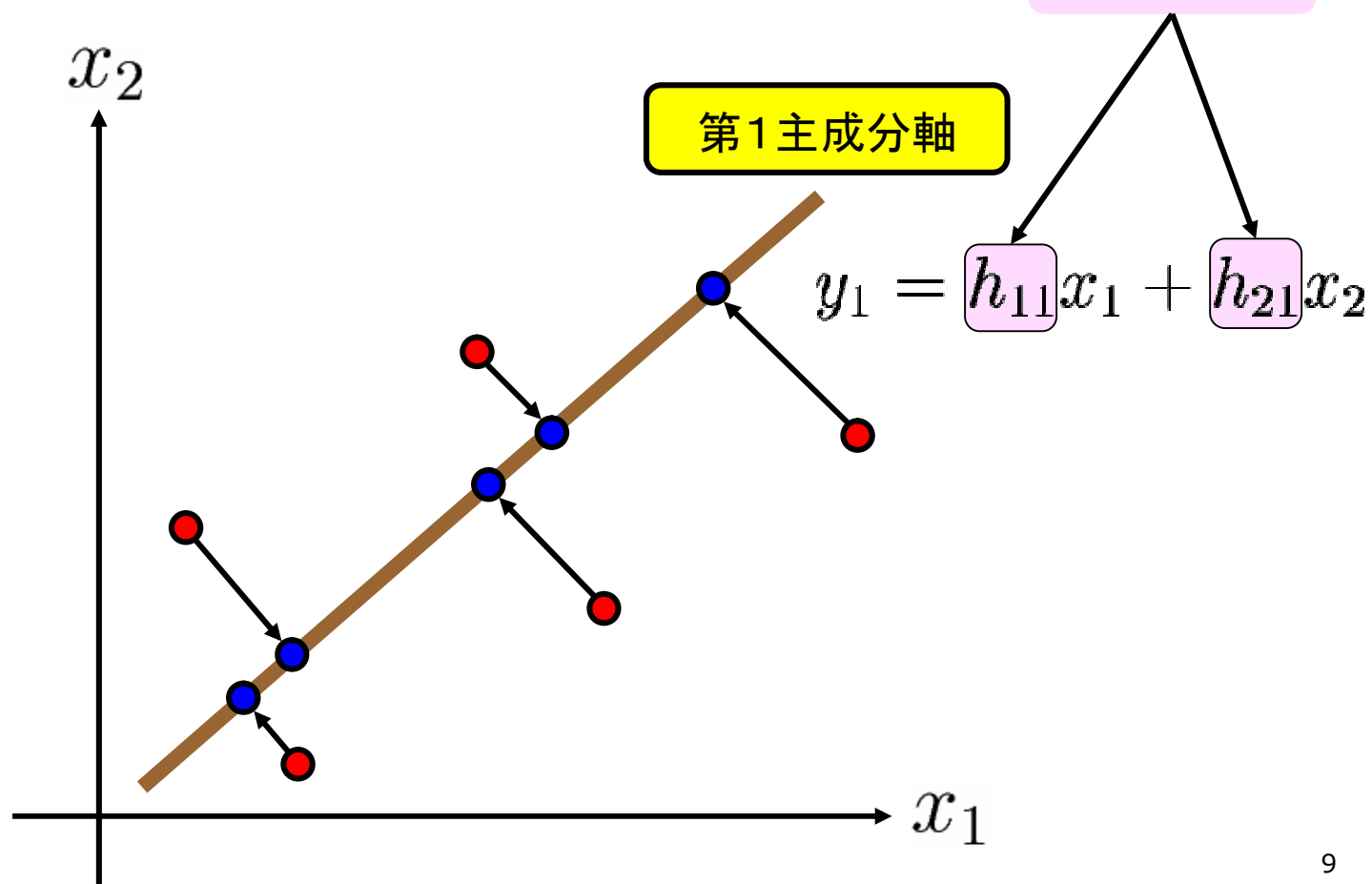


情報の損失が多い

射影したデータの分散が最大となる軸を探す

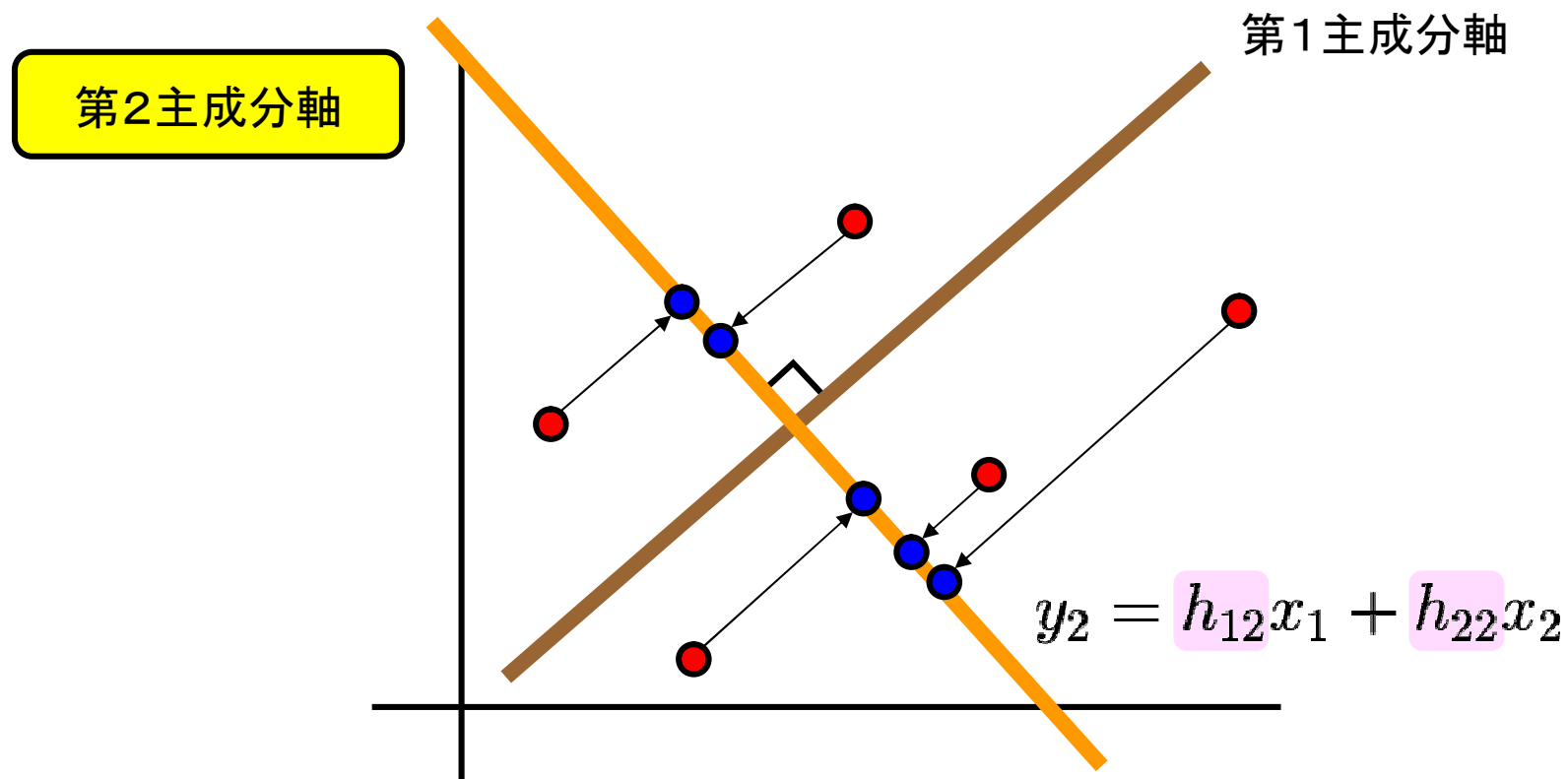
第1主成分

- 射影したデータの分散が最大となるような軸を探す



第2主成分

- 第1主成分と直交する軸の中で、軸上に射影したデータの分散が最大となる軸を探す

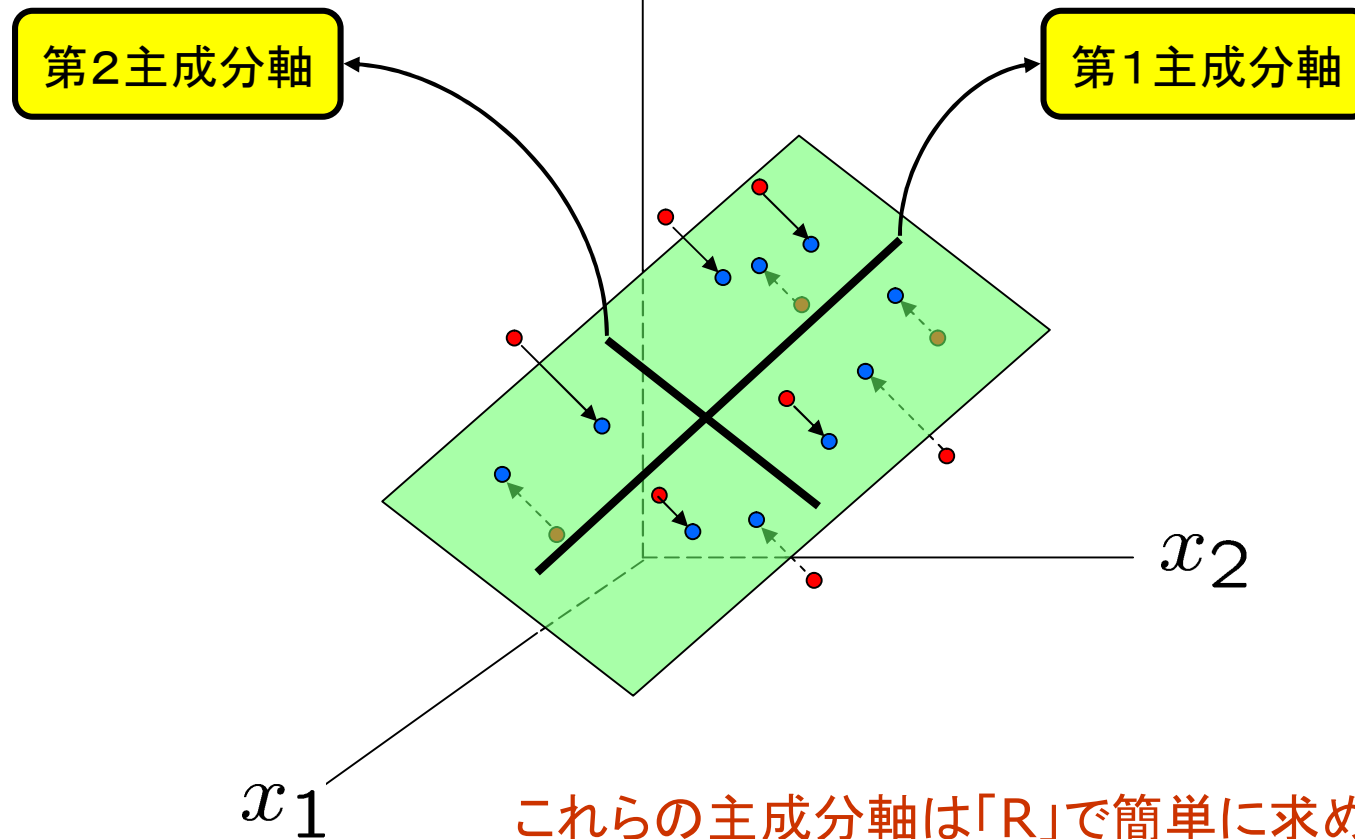


3次元空間から2次元空間への縮約

$$y_2 = h_{12}x_1 + h_{22}x_2 + h_{32}x_3$$

x_3

$$y_1 = h_{11}x_1 + h_{21}x_2 + h_{31}x_3$$



主成分分析の要点

- 主成分の分散：主成分がもつ情報量
＝分散共分散行列(相関行列)の固有値
- 主成分を構成する係数は、分散共分散行列(相関行列)の固有ベクトルを求める手続きにより得られる
- 主成分分析では、分散共分散行列から分析を行う場合と、相関行列から行う場合で結果が異なる。
- データが異なる尺度(単位)で測定されている場合には、変数を基準化して分析を行う必要がある。

具体例：成績データ

杉山 高一 著「多変量データ解析入門」

□ 中学2年生の成績データ

- 標本数：166
- 変数の数：科目数=9
 - 国語、社会、数学、理科、音楽、英語、体育、技家、英語
- ダウンロードしたファイルに記載されている最後の3列のデータ(変数名:「4year」「5year」「6year」)を削除して分析を行う
- 9科目の得点を適当に組み合わせた変数を作り、できるだけ少ない変数で生徒の特徴を捉えたい

データのダウンロード

□ 統計科学研究所のウェブサイト

- <http://www.statistics.co.jp/index.htm>



統計科学研究所

研究テーマ | 統計データ分析士資格認定 |
セミナー開催の趣旨 | セミナー | 開講科目詳細 | セミナー申し込み | 問い合わせ・資料請求 |

統計科学研究所

- 研究テーマ
- 統計データ分析士資格認定
- 特別セミナー開催の趣旨
- 特別セミナー開講科目
- 開講科目詳細と講師紹介
- セミナー申し込み
- 問い合わせ・資料請求
- 個人情報保護方針
- 統計分析ソフトウェア「R」
- データ

データのダウンロード

- 成績のデータの[csv]を右クリック
⇒名前を付けて保存



統計科学研究所

研究テーマ | 統計データ分析士資格認定 |
セミナー開催の趣旨 | セミナー | 開講科目詳細 | セミナー申し込み | お問い合わせ・資料請求 | Top |

データ

1. 男物着尺地と紳士服地のデータ [xls] [csv]
2. 京の字のデータ [xls] [csv]
3. Aboriginesの手のデータ [xls] [csv]
4. 白人の手のデータ [xls] [csv]
5. 成績のデータ [xls] [csv]

Copyright © Toukei Kagaku Kenkyujo, Co., Ltd. All right reserved.

主成分分析を行うプログラム

```
seiseki <- read.csv("seiseki.csv", header=T)
result <- prcomp(seiseki, scale=T)
summary(result)
biplot(result)
```

□ プログラムの概要

- 1行目: データの読み込み
- 2行目: 主成分分析を行う関数 **"prcomp"** を適用
- 3行目: 主成分分析の結果の要約の出力
- 4行目: 主成分得点をプロットする関数 **"biplot"** を適用

相関行列から主成分分析を行う

```
result <- prcomp(seiseki, scale=T)
```

□ 引数 "scale" について

- 関数 "prcomp" に、引数 "scale=T" を指定
⇒ 相関行列から主成分分析を行う
- 関数 "prcomp" に、引数 "scale=F" を指定
⇒ 分散共分散行列から主成分分析を行う

分析結果の要約

- 分析結果に関数 “summary” を適用
 - Standard deviation (標準偏差)
 - Proportion of Variance (寄与率)
 - Cumulative Proportion (累積寄与率)

```
R Console
> summary(result)
Importance of components:

```

	PC1	PC2	PC3	PC4	PC5	PC6	PC7	PC8	PC9
Standard deviation	2.451	1.048	0.7006	0.6379	0.5480	0.4706	0.4275	0.4138	0.3491
Proportion of Variance	0.667	0.122	0.0545	0.0452	0.0334	0.0246	0.0203	0.0190	0.0135
Cumulative Proportion	0.667	0.789	0.8439	0.8892	0.9225	0.9471	0.9674	0.9865	1.0000

```
> |
```

第1主成分には、全体の67%の情報が縮約されている

第4主成分までで、全体の約90%の情報を占める

寄与率と累積寄与率

□ 標準偏差 該当する主成分がもつ情報量

- 第 i 主成分の標準偏差。これらは第 i 固有値の平方根となっている。

□ 第 i 主成分の寄与率

- $\frac{l_i}{\text{総分散}}$, 総分散 = $\sum_{i=1}^p l_i$, l_i : 第 i 固有値

全情報量のうち、該当する主成分が占める情報量の割合

□ 第 m 主成分までの累積寄与率

- $\frac{l_1 + l_2 + \dots + l_m}{\text{総分散}}$ 選択した主成分が占める情報量の割合
次元の縮約により失う情報量を測ることができる

分析結果の出力

- 次のようにして、関数 “prcomp” で得られたオブジェクトから、分析結果を得ることができる
- 今回のプログラムの場合
 - `result$rotation` : 固有ベクトル(主成分軸の係数)
 - `result$x` : 主成分得点
- 関数 “round” を使って出力結果を適当な桁数で丸めると見やすくなる
 - `round(result$x, digits=3)` : 主成分得点を小数点3桁で表示

固有ベクトルの出力

- `round(result$rotation, 3)` の出力

PC : Principal Component ↔ 主成分

```
R Console
> round(result$rotation, 3)
      PC1  PC2  PC3  PC4  PC5  PC6  PC7  PC8  PC9
kokugo 0.363 -0.149 0.074 -0.236 0.301 -0.494 0.620 0.110 -0.231
shakai 0.369 0.147 -0.062 -0.107 0.087 -0.573 -0.517 -0.235 0.412
sugaku 0.357 0.181 -0.400 0.029 0.061 0.408 0.409 -0.446 0.377
rika    0.367 0.251 0.008 0.067 -0.262 0.039 -0.177 -0.392 -0.736
ongaku  0.354 -0.010 -0.200 0.357 -0.642 -0.132 0.119 0.495 0.133
bijutu  0.313 -0.312 0.264 0.712 0.440 0.136 -0.125 0.002 -0.003
taiiku  0.139 -0.859 -0.080 -0.284 -0.269 0.107 -0.128 -0.235 0.007
gika    0.317 0.149 0.784 -0.293 -0.190 0.287 0.042 0.062 0.231
eigo    0.357 0.047 -0.317 -0.355 0.338 0.361 -0.320 0.525 -0.146
> |
```

第1主成分の構成

□ 第1主成分 =

$$0.363 \times \text{国語} + 0.369 \times \text{社会} + 0.357 \times \text{数学} + \\ 0.367 \times \text{理科} + 0.354 \times \text{音楽} + 0.313 \times \text{美術} + \\ 0.139 \times \text{体育} + 0.317 \times \text{技家} + 0.357 \times \text{英語}$$

```
R Console
> round(result$rotation, 3)
      PC1    PC2    PC3    PC4    PC5    PC6    PC7    PC8    PC9
kokugo 0.363 -0.149  0.074 -0.236  0.301 -0.494  0.620  0.110 -0.231
shakai 0.369  0.147 -0.062 -0.107  0.087 -0.573 -0.517 -0.235  0.412
sugaku 0.357  0.181 -0.400  0.029  0.061  0.408  0.409 -0.446  0.377
rika   0.367  0.251  0.008  0.067 -0.262  0.039 -0.177 -0.392 -0.736
ongaku 0.354 -0.010 -0.200  0.357 -0.642 -0.132  0.119  0.495  0.133
bijutu 0.313 -0.312  0.264  0.712  0.440  0.136 -0.125  0.002 -0.003
taiiku 0.139 -0.859 -0.080 -0.284 -0.269  0.107 -0.128 -0.235  0.007
gika   0.317  0.149  0.784 -0.293 -0.190  0.287  0.042  0.062  0.231
eigo   0.357  0.047 -0.317 -0.355  0.338  0.361 -0.320  0.525 -0.146
> |
```

第1主成分の構成

□ 第1主成分 =

$$0.363 \times \text{国語} + 0.369 \times \text{社会} + 0.357 \times \text{数学} + \\ 0.367 \times \text{理科} + 0.354 \times \text{音楽} + 0.313 \times \text{美術} + \\ \text{小 } 0.139 \times \text{体育} + 0.317 \times \text{技家} + 0.357 \times \text{英語}$$

```
R Console
> round(result$rotation, 3)
      PC1      PC2      PC3      PC4      PC5      PC6      PC7      PC8      PC9
kokugo 0.363 -0.149  0.008  0.067 -0.262  0.039 -0.177 -0.392 -0.736
shakai 0.369  0.147 -0.200  0.357 -0.642 -0.132  0.119  0.495  0.133
sugaku 0.357  0.181  0.000  0.000  0.000  0.000  0.000  0.000 -0.003
rika   0.367  0.251  0.000  0.000  0.000  0.000  0.000  0.000  0.007
ongaku 0.354 -0.010  0.000  0.000  0.000  0.000  0.000  0.000  0.231
bijutu 0.313 -0.312  0.000  0.000  0.000  0.000  0.000  0.000 -0.146
taiiku 0.139 -0.859  0.000  0.000  0.000  0.000  0.000  0.000  0.000
gika   0.317  0.149  0.000  0.000  0.000  0.000  0.000  0.000  0.231
eigo   0.357  0.047  0.000  0.000  0.000  0.000  0.000  0.000 -0.146
> |
```

第1主成分
⇒ 筆記試験の総合得点の因子

第1主成分が大きい
⇒ 筆記試験の総合得点が高い

第2主成分の解釈

□ 第2主成分 =

$$\begin{aligned} & -0.149 \times \text{国語} + 0.147 \times \text{社会} + 0.181 \times \text{数学} \\ & + 0.251 \times \text{理科} - 0.010 \times \text{音楽} - 0.312 \times \text{美術} \\ & - 0.859 \times \text{体育} + 0.149 \times \text{技家} + 0.047 \times \text{英語} \end{aligned}$$

```
R Console
> round(result$rotation, 3)
      PC1  PC2  PC3  PC7  PC8  PC9
kokugo 0.363 -0.149 0.074 0.620 0.110 -0.231
shakai 0.369 0.147 -0.062 0.517 -0.235 0.412
sugaku 0.357 0.181 -0.400 0.409 -0.446 0.377
rika    0.367 0.251 0.008 0.067 -0.262 0.039 -0.177 -0.392 -0.736
ongaku 0.354 -0.010 -0.200 0.357 -0.642 -0.132 0.119 0.495 0.133
bijutu 0.313 -0.312 0.264
taiiku 0.139 -0.859 -0.080
gika    0.317 0.149 0.784
eigo    0.357 0.047 -0.317
> |
```

第2主成分
⇒ 体育の因子

第2主成分が小さい（符号に注意）
⇒ 体育の得点が優れている

因子負荷量

□ 各主成分の意味づけ

- 主成分に強く寄与している変数を見つけることが重要

□ 因子負荷量

- 主成分と各変数との相関係数

相関行列から分析を始めた場合の因子負荷量

$$\text{Cor}(x_i, y_j) = \sqrt{l_j} h_{ij}, \quad x_i : i \text{ 番目の変数}, \quad y_j : \text{第 } j \text{ 主成分}$$

参考：奥野 忠一著「多変量解析法 改訂版」日科技連

- 因子負荷量が1か-1に近い因子ほど、主成分に強く寄与している
- 因子負荷量をプロットすることにより、主成分に寄与している因子を視覚的に捉えることができる

因子負荷量に関するプログラム

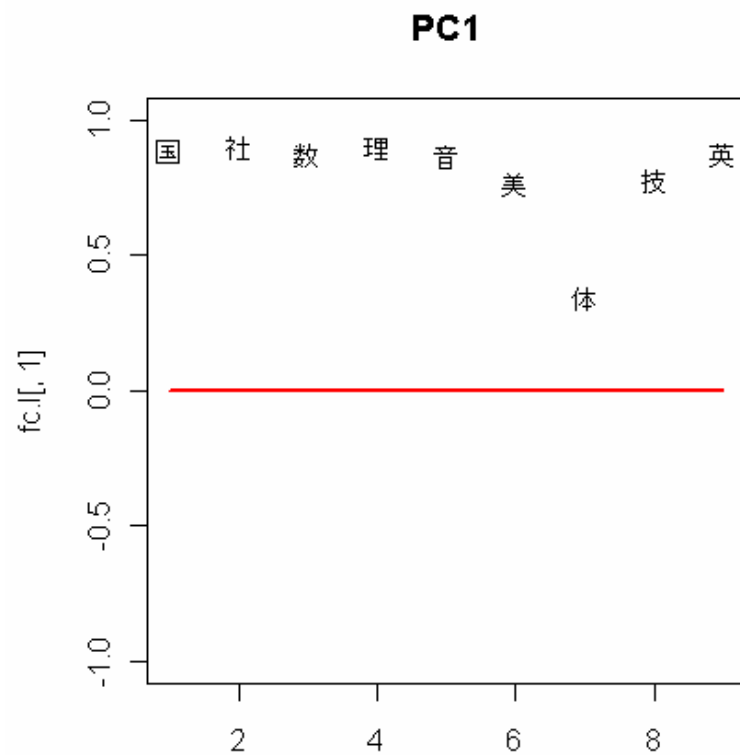
```
fc.l <- sweep(result$rotation, MARGIN=2, result$sdev, FUN="*")
subject <- c("国", "社", "数", "理", "音", "美", "体", "技", "英")
plot(fc.l[,1], pch=subject, ylim=c(-1,1), main="PC1")
plot(fc.l[,2], pch=subject, ylim=c(-1,1), main="PC2")
```

□ プログラムの概略

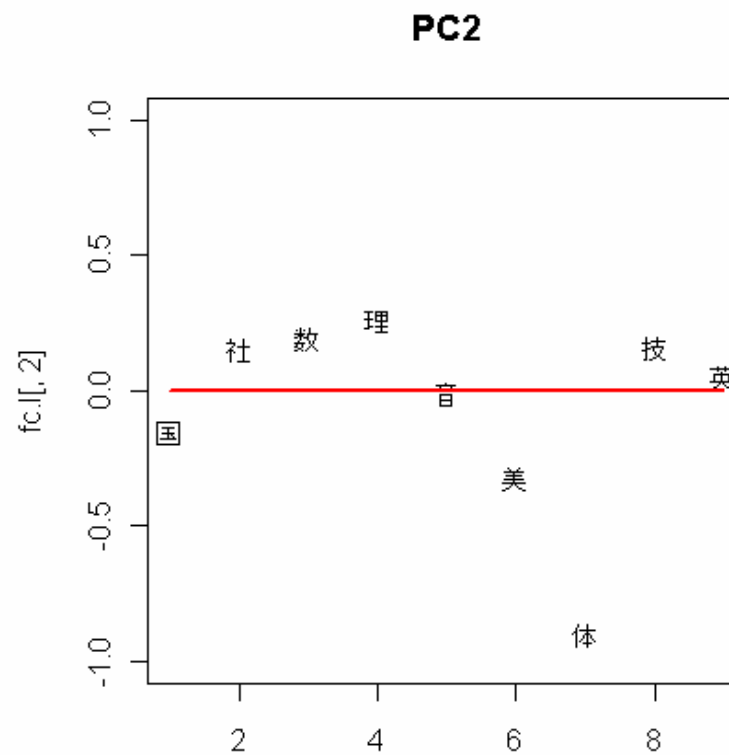
- 1行目：因子負荷量の計算
 - 固有ベクトル(result\$rotation)と、対応した固有値の平方根(result\$sdev)との積をとる
- sweep 関数の使い方は、apply 関数とよく似ている
 - 参考URL：R-Tips 24節「applyファミリー」

<http://cse.naro.affrc.go.jp/takezawa/r-tips/r/24.html>

因子負荷量のプロット（1次元）

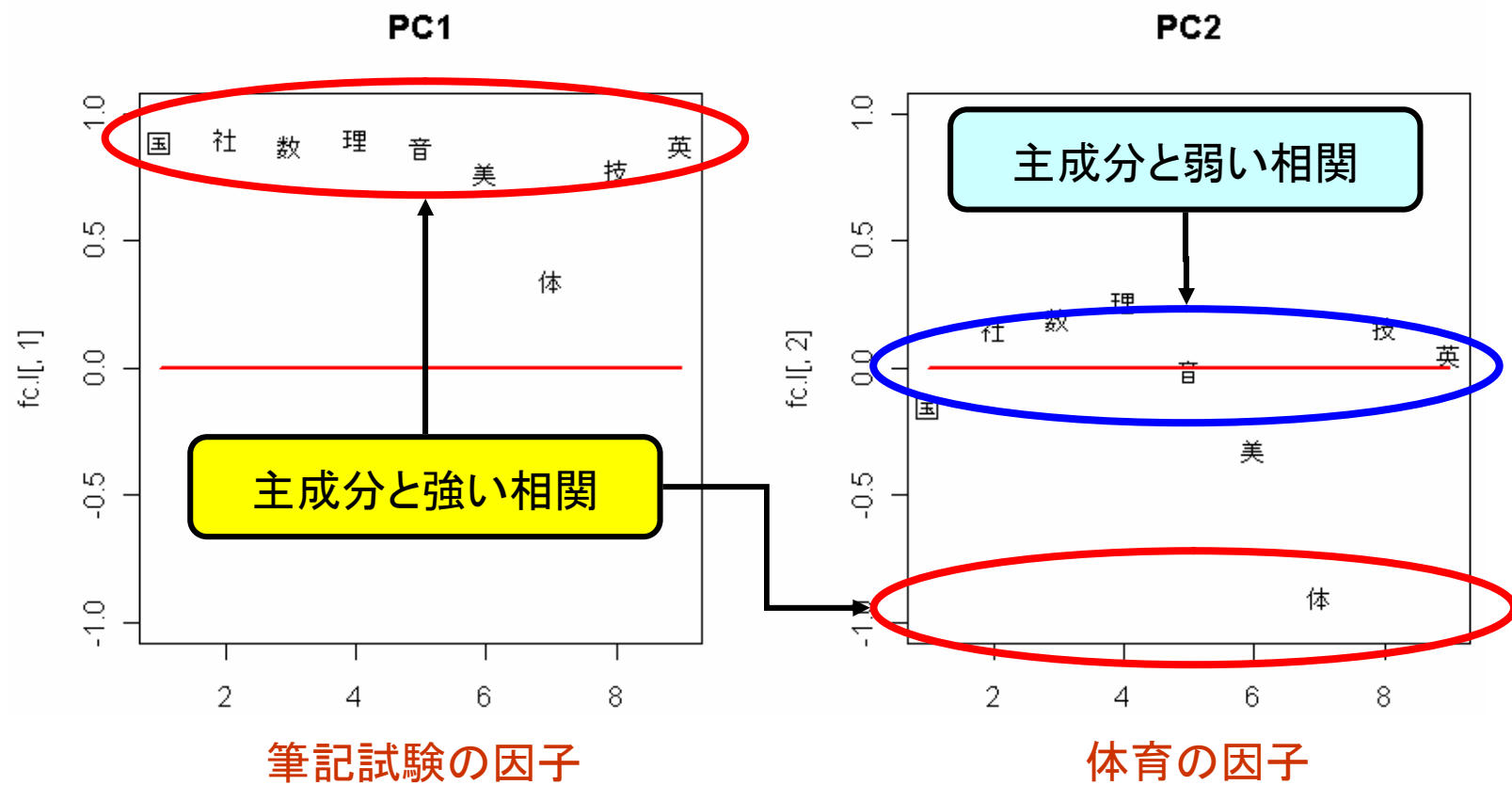


3行目のプログラムの出力



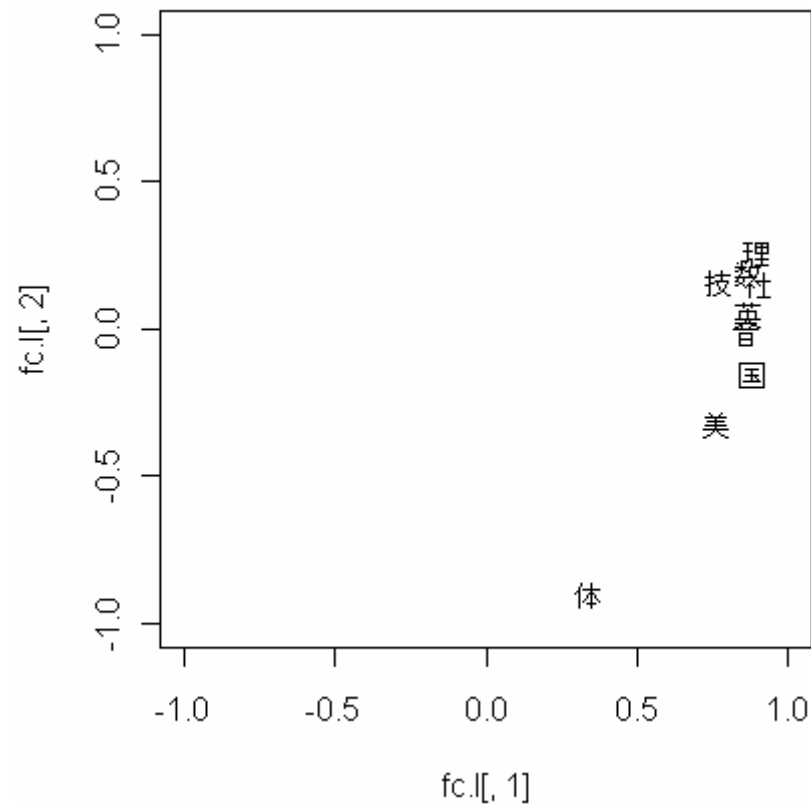
4行目のプログラムの出力

因子負荷量の解釈（1次元）

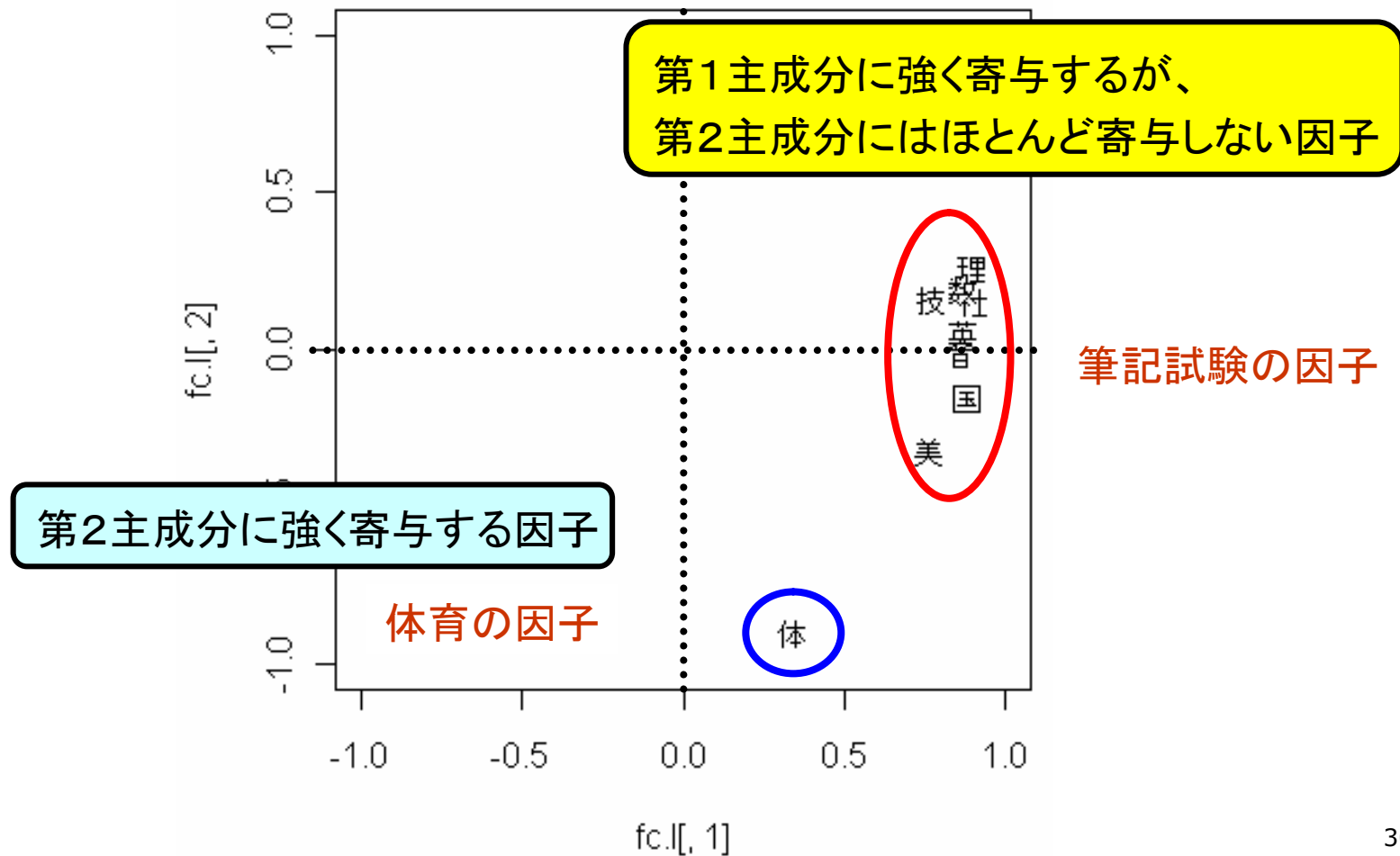


因子負荷量のプロット（2次元）

```
plot(fc.l[,1], fc.l[,2], pch=subject,  
     xlim=c(-1,1), ylim=c(-1,1), main=title)
```



因子負荷量の解釈（2次元）



主成分得点の定義

□ 主成分得点の定義

- 主成分 $y_i = h_{i1}x_1 + h_{i2}x_2 + \dots + h_{ip}x_p$ に個々のデータを代入したもの

result\$x

□ 成績データの例

- 第1主成分 =

$$0.363 \times \text{国語} + 0.369 \times \text{社会} + \boxed{0.357 \times \text{数学}} + \\ 0.367 \times \text{理科} + 0.354 \times \text{音楽} + 0.313 \times \text{美術} + \\ 0.139 \times \text{体育} + 0.317 \times \text{技家} + \boxed{h_{13}x_3} \times \text{英語}$$

国語	社会	数学	理科	音楽	美術	体育	技家	英語
95	87	77	100	77	82	78	96	87

相関行列から分析を行う場合は、全ての変数を基準化したものを代入する 31

主成分得点

□ 4人目の成績

国語	社会	数学	理科	音楽	美術	体育	技家	英語
95	87	77	100	77	82	78	96	87

第1主成分得点 : 5.107

第2主成分得点 : 0.228

□ 130人目の成績

国語	社会	数学	理科	音楽	美術	体育	技家	英語
64	36	20	31	53	68	99	7	26

第1主成分得点 : -0.812

第2主成分得点 : -2.244

主成分得点の出力

- result\$x : 主成分得点を出力する

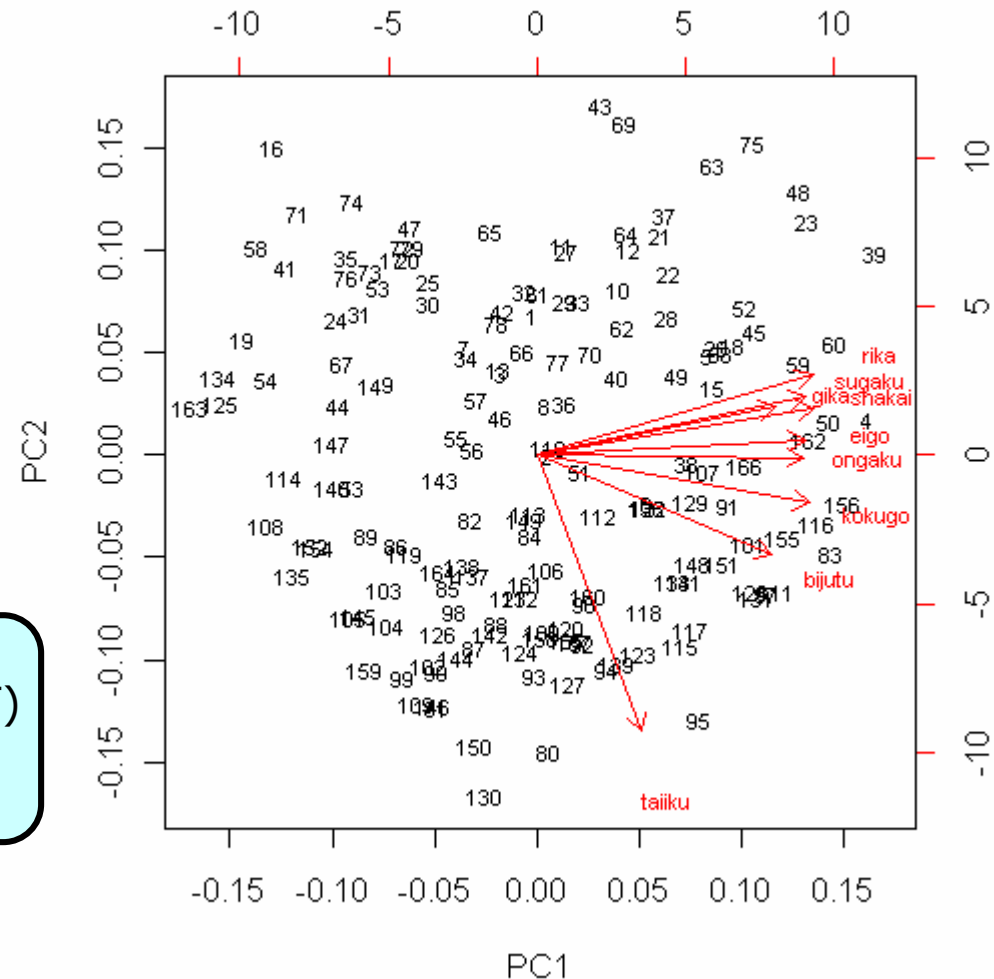
```
R Console
> round(result$x, 3)
      PC1    PC2    PC3    PC4    PC5    PC6    PC7    PC8    PC9
[1,] -0.056  0.917 -0.116  1.189 -0.828  0.097 -0.367 -0.272  0.091
[2,]  0.148  0.000  0.660  1.178 -1.237  0.693  0.439  0.209  0.060
[3,] -0.540  0.543  0.963  1.847 -1.092 -0.049 -0.341  0.500 -0.156
[4,]  5.107  0.228  0.626 -0.580 -0.293 -0.369 -0.302 -0.243 -0.313
[5,]  2.637  0.661  0.012  0.880 -0.047 -0.465  0.195 -0.287  0.632
[6,]  1.691 -0.326  0.493  0.288 -0.301 -0.196  0.142 -0.104  0.177
[7,] -1.096  0.711  0.325  1.106 -0.128  0.469 -0.054 -0.349 -0.065
[8,]  0.116  0.324  1.249  0.036 -0.265 -0.702 -0.106 -0.781  0.106
[9,] -2.243  0.107 -0.483  0.078 -0.020 -0.670
[10,]  1.244 -0.119 -0.276 -0.067  0.607 -1.210
```

例で見た4番目の生徒の
主成分得点

主成分得点のプロット (biplot)

- 主成分得点を低次元空間にプロットすると、個体の特徴や位置を把握しやすくなる
- 「R」では、biplot 関数を適用することで、解釈しやすい形で主成分得点のプロットを得ることができる

```
result <- prcomp(seiseki, scale=T)  
biplot(result)
```



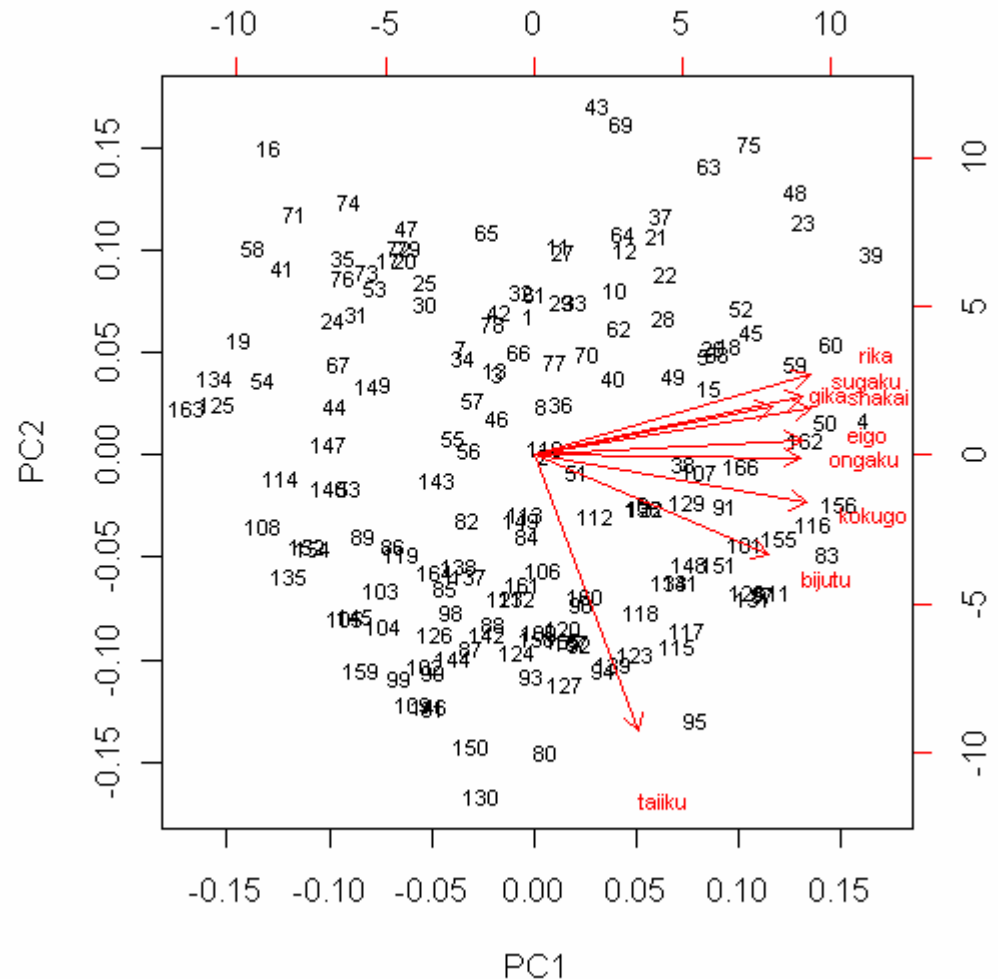
biplotの解釈

□ 第一主成分の解釈

- 筆記試験の総合得点
⇒右にあるデータほど筆記試験の総合点が高い

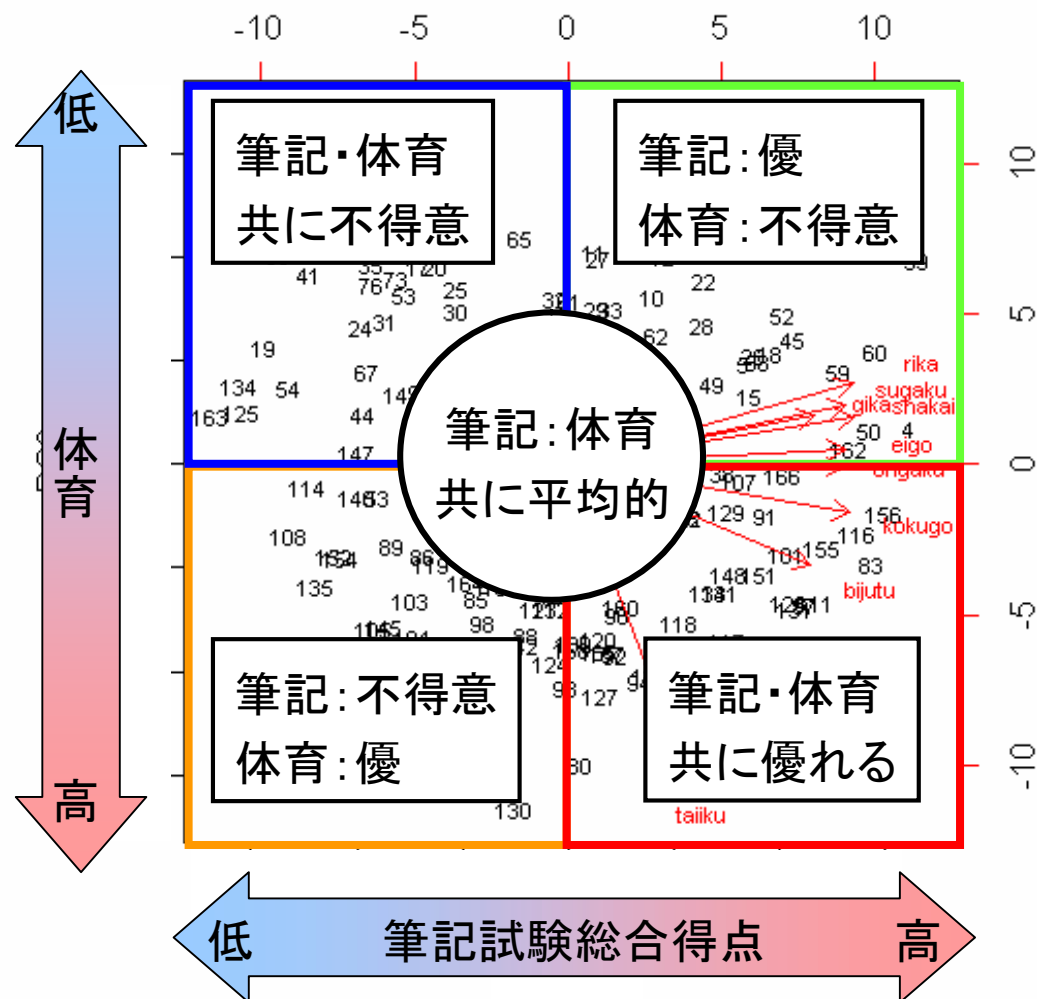
□ 第二主成分の解釈

- 体育の因子
⇒下にあるデータほど体育の成績が良い
- 主成分の符号やベクトルの向きに注意す。



主成分によるデータの位置づけ

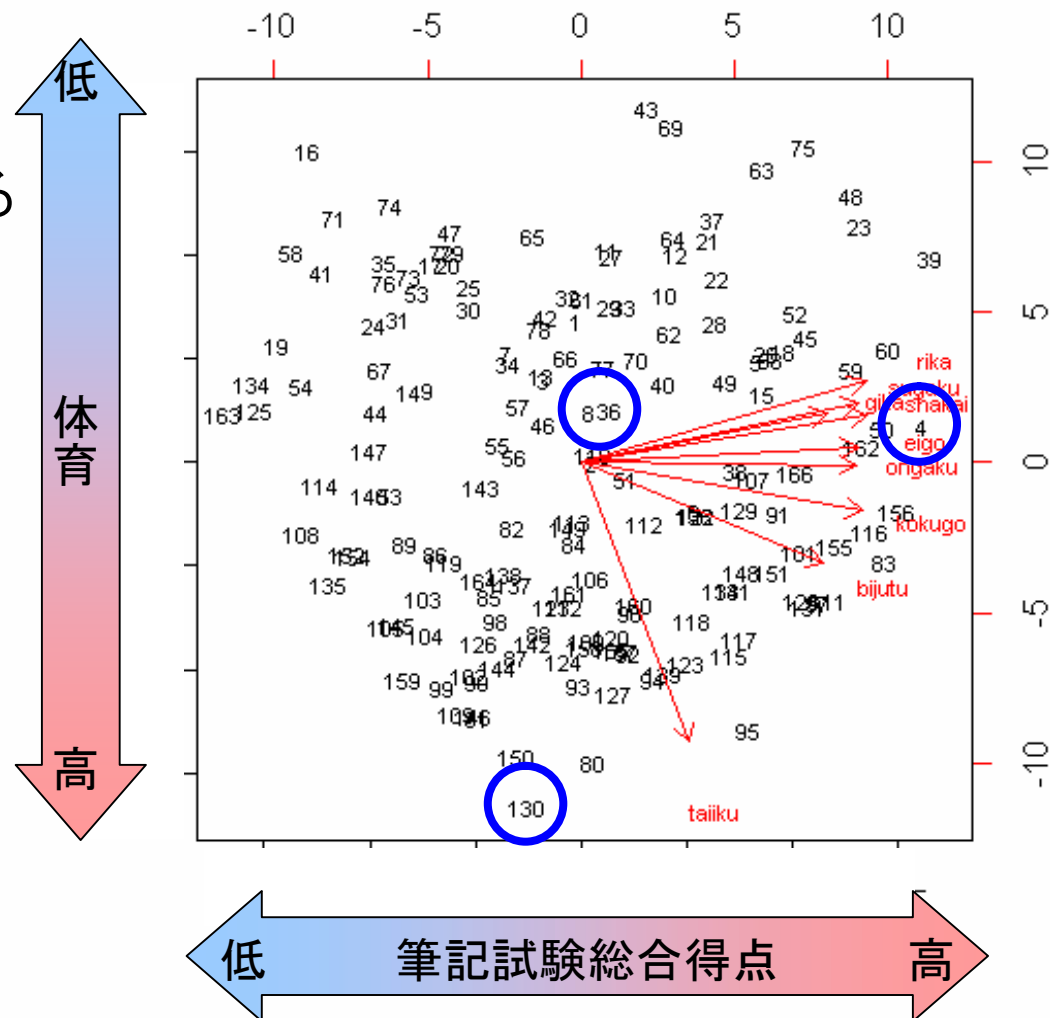
- 主成分の解釈から、各個体を右図のように分類して考えることができる
- 主成分得点の順にデータを並べ替えることである特性値について順位付けを行うこともできる



主成分によるデータの位置づけ

□ biplotの見方

- 4番
筆記試験が優れている
体育は平均程度
- 130番
筆記試験は平均程度
体育得意
- 8番
筆記試験も体育も
平均程度



主成分得点とデータ

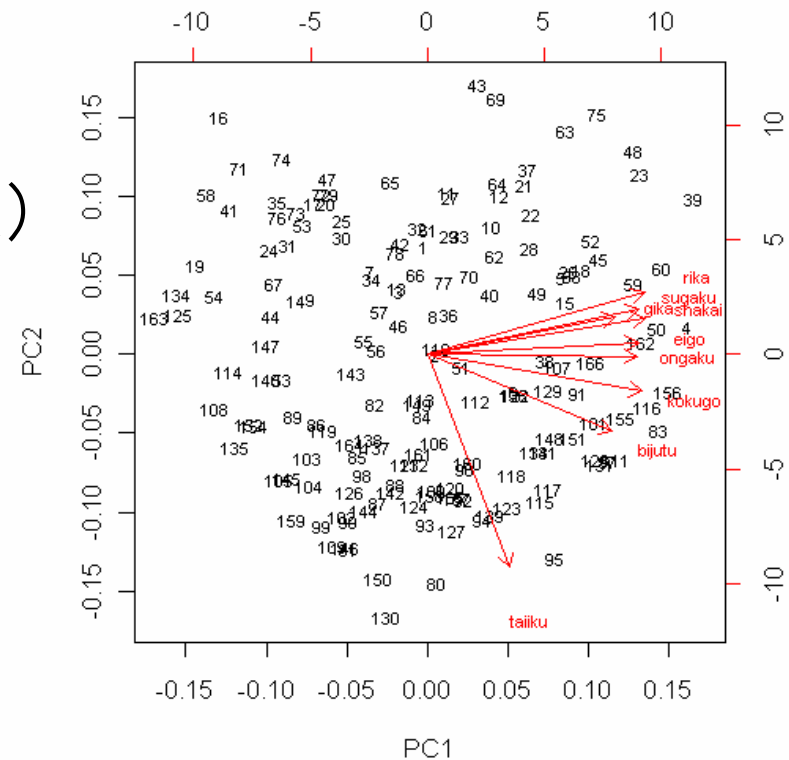
□ 例で挙げた生徒の成績と主成分得点

- PC1 : 第1主成分の主成分得点
- PC2 : 第2主成分の主成分得点

	国語	社会	数学	理科	音楽	美術	体育	技家	英語	PC1	PC2
4	95	87	77	100	77	82	78	96	87	5.1	0.2
8	56	54	37	59	35	64	53	67	7	0.1	0.3
130	64	36	20	31	53	68	99	7	26	-0.8	-2.2
平均	57.5	39.6	45.6	49.9	42.6	62.5	57.7	47.3	39.1	0	0

まとめ

- 主成分分析を行う関数 “**prcomp**” の使い方
- 主成分得点の出力の仕方
 - **obj\$x**
- 固有ベクトル(主成分軸の係数)の出力の仕方
 - **obj\$rotation**
- 因子負荷量の求め方と解釈
- 関数 “**biplot**”の使い方と解釈



参考URL

- 統計科学研究所のウェブサイト

<http://www.statistics.co.jp/index.htm>

- R-Tips

<http://cse.naro.affrc.go.jp/takezawa/r-tips/r2.html>

- JIN'S PAGE

<http://www1.doshisha.ac.jp/~mjn/R/>