

# 主成分分析

1. データ  $X$  の主成分分析は、各点からの垂線の距離（情報損失量）が最小になるように、固有ベクトルを決める。
2. 主成分得点の分散の最大化も、これと同じことをしている。
3. 固有値は、最小化された情報損失量であり、最大化された主成分得点の分散でもある。
4. Rで主成分分析は `prcomp()` 関数を使う。
5. 固有値の平方根（主成分得点の標準偏差） `$sdev` とノルム 1 で直交する固有ベクトル `$rotation` が出力される。
6. 主成分得点（スコア）は `$x` で出力される。
7. 通常、主成分分析では、各変数の分散は1に標準化した方がよいが、これは引数 `scale=1` ができる。
8. 主成分負荷量は、標準化した変数と、それを使った主成分分析のスコアとの相関係数とするのが一般的。
9. `biplot()` 関数を使うと、スコアと主成分負荷量という別種の値が1つのグラフにプロットされるが、これらの値は、特異値分解で得られるもので、上記のそれらとは値が異なる。
10. `biplot()` 関数によるプロット図は、引数 `scale=` に0~1の値で変わる。スコアを一般的な値で出力するなら `scale=0`。負荷量のプロットを  $X$  とスコアの相関係数の比率（値そのものは異なる）でプロットしたいなら `scale=1`（これがデフォルト）

## 考え方

主成分分析のかみ砕いた解説（咀嚼説明）は、有馬・石村（1987）<sup>1</sup> にあります。

以下はそのフォローです。

成分ベクトルと  $X$  との距離の最小化であり、

スコアの分散の最大化でもある。

## 情報損失量 $U$ の最小化

$x_i = (x_{i1}, x_{i2})$  と直線  $a_2 x_1 - a_1 x_2 + a_0 = 0$  との距離を  $l$  で表す。

ただし、 $a_1^2 + a_2^2 = 1$  とする。

ヘッセの標準形から

$$l_i = |a_2 x_{i1} - a_1 x_{i2} + a_0|$$

ヘッセの標準形

点  $(x_1, y_1)$  から直線  $ax + by + c = 0$  に下した垂線の長さは

$$\frac{|ax_1 + by_1 + c|}{\sqrt{a^2 + b^2}}$$

情報損失量  $U$  を以下で定義して、

$$U(a_2, a_1, a_0) = \sum_i l_i^2 = \sum_i (a_2 x_{i1} - a_1 x_{i2} + a_0)^2$$

$U$  を最小化するような  $a_2, a_1, a_0$  を見つけたい。

言い換えれば・・・

$$\min_{a_2, a_1, a_0} U(a_2, a_1, a_0)$$

$$\text{s.t. } a_2^2 + a_1^2 = 1$$

なので、ラグランジュ乗数法より

$$\min_{a_2, a_1, a_0} L = U - \lambda_0 (a_2^2 + a_1^2 - 1) \quad (1)$$

1階の条件として

$$\frac{\partial U}{\partial a_0} = 2 \sum_i (a_2 x_{i1} - a_1 x_{i2} + a_0) = 0 \quad (2)$$

なので、 $x_{i1}$  と  $x_{i2}$  の平均を、それぞれ  $\bar{X}_1, \bar{X}_2$  とすると、

$$a_2 \bar{X}_1 - a_1 \bar{X}_2 + a_0 = 0$$

ということであり、 $z_{i1} = x_{i1} - \bar{X}_1, z_{i2} = x_{i2} - \bar{X}_2$  とすると、

$$a_2 x_{i1} - a_1 x_{i2} + a_0 = a_2 z_{i1} - a_1 z_{i2}$$

なので、

$$\frac{\partial U}{\partial a_2} = 2 \sum_i z_{i1} (a_2 z_{i1} - a_1 z_{i2}) - 2\lambda_0 a_2 = 0 \quad (3)$$

$$\frac{\partial U}{\partial a_1} = -2 \sum_i z_{i2} (a_2 z_{i1} - a_1 z_{i2}) - 2\lambda_0 a_1 = 0 \quad (4)$$

さらに、

$$\frac{\partial U}{\partial \lambda_0} = -(a_2^2 + a_1^2 - 1) = 0 \quad (5)$$

(3)式から

$$a_2 \text{Var}(X_1) - a_1 \text{Cov}(X_1, X_2) = \frac{\lambda_0}{n-1} a_2 \quad (6)$$

(4)式から

$$-a_2 \text{Cov}(X_1, X_2) + a_1 \text{Var}(X_2) = \frac{\lambda_0}{n-1} a_1 \quad (7)$$

$\lambda = \lambda_0 / (n-1)$  と書き換えると、

$$\begin{bmatrix} \text{Var}(X_1) & -\text{Cov}(X_1, X_2) \\ -\text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{bmatrix} \begin{bmatrix} a_2 \\ a_1 \end{bmatrix} = \lambda \begin{bmatrix} a_2 \\ a_1 \end{bmatrix} \quad (8)$$

ということであり、 $(a_2, a_1)$ というのは、 $n$ 行2列の行列 $X = (X_1, -X_2)$ の分散・共分散行列の**固有ベクトル**！

$\lambda$ は、 $X$ の分散・共分散行列の**固有値**！

・・・ということ。

さらに、(6)式に $a_2$ を掛けて、(7)式に $a_1$ を掛けて足し合わせると

$$a_2^2 \text{Var}(X_1) - 2a_1 a_2 \text{Cov}(X_1, X_1) + a_1^2 \text{Var}(X_2) = \lambda(a_2^2 + a_1^2)$$

なので、

$$\text{Var}(a_2 X_1 - a_1 X_2) = \lambda \quad (9)$$

であり、これは、

$$\frac{1}{n-1} \sum_i (a_2 x_i - a_1 x_{2i} + a_0)^2 = \lambda$$

のことなので、 $\lambda_0$ は**情報損失量** $U$ ということ。つまり、

$$U = \lambda_0 \quad (10)$$

## スコアの最大化

$a_1 X_1 + a_2 X_2$ を**主成分得点** (principal component score)

あるいは単に**スコア** (score) と呼ぶ。

$n$ 行2列の行列 $X = (X_1, X_2)$ でスコア $a_1 X_1 + a_2 X_2$ を最大化するノルム1の $(a_1, a_2)$ を求める。

$$\max_{a_1, a_2} \text{Var}(a_1 X_1 + a_2 X_2)$$

$$1. \quad 1. \quad a_1^2 + a_2^2 = 1$$

であるから、ラグランジュ乗数法で

$$\max_{a_1, a_2} L = \text{Var}(a_1 X_1 + a_2 X_2) - \mu(a_1^2 + a_2^2 - 1)$$

$$= a_1^2 \text{Var}(X_1) + 2a_1 a_2 \text{Cov}(X_1, X_2) + a_2^2 \text{Var}(X_2) - \mu(a_1^2 + a_2^2 - 1) \quad (10)$$

1階の条件は

$$\frac{\partial L}{\partial a_1} = 2a_1 \text{Var}(X_1) + 2a_2 \text{Cov}(X_1, X_2) - 2\mu a_1 = 0 \quad (11)$$

$$\frac{\partial L}{\partial a_2} = 2a_1 \text{Cov}(X_1, X_2) + 2a_2 \text{Var}(X_2) - 2\mu a_2 = 0$$

(12)

したがって,

$$\begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \mu \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

(13)

であり,  $(a_1, a_2)$ は,  $X = (X_1, X_2)$  の分散・共分散行列の固有ベクトル

$\mu$ は, その固有値.

$$a_1^2 \text{Var}(X_1) + 2a_1 a_2 \text{Cov}(X_1, X_2) + a_2^2 \text{Var}(X_2) = \mu(a_1^2 + a_2^2)$$

なので

$$\text{Var}(a_1 X_1 + a_2 X_2) = \mu \quad (14)$$

でもある.

## Rで

```
1 library(tidyverse)
```

Rの `pcrcom` 関数で

## データ

わかりやすいデータを作成.

$X = (X_1, X_2)$ の2変数9個のデータ.

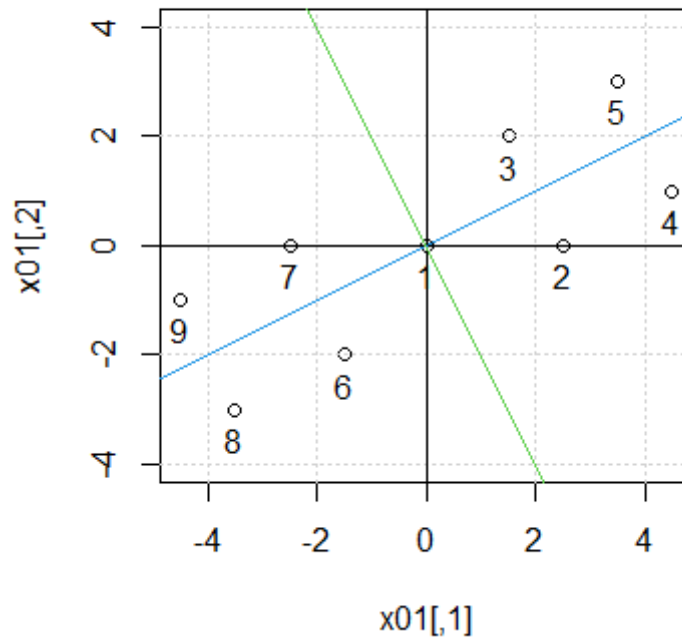
```
1 x01<-c(0,0,
2 2+0.5,1-1,
3 2-0.5,1+1,
4 4+0.5,2-1,
5 4-0.5,2+1,
6 -2+0.5,-1-1,
7 -2-0.5,-1+1,
8 -4+0.5,-2-1,
9 -4-0.5,-2+1
10 )%>%
11 matrix(ncol=2,byrow=T)
```

こんな感じ

```

1 plot(x01,asp=1)
2 grid()
3 text(1:9,x=x01[,1],y=x01[,2],pos = 1)
4 abline(h=0)
5 abline(v=0)
6 abline(a=0,b=0.5,col=4)
7 abline(a=0,b=-2,col=3)

```



横軸が  $X_1$  (`x01[,1]`) , 縦軸が  $X_2$  (`x01[,2]`)

8つの点が,  $(2, 1)$  の方向の原点を通る直線に対して垂直に, 長さ  $\sqrt{0.5^2 + 1} = \sqrt{1.25}$  で対称にばらついている。

分析しなくても主成分は,  $(2, 1)$  の方向と  $(-1, 2)$  の方向であることはわかる。

平均は2変数とも0。

```

1 apply(x01,2,mean)%>%
2 round(4)

```

```
1 ## [1] 0 0
```

分散・共分散行列  $M$  は

```
1 var(x01)
```

```

1 ##      [,1] [,2]
2 ## [1,] 10.25 4.5
3 ## [2,] 4.50 3.5

```

# 主成分分析を行う

`prcomp()` 関数

```
1 pr01<-prcomp(x01)
2 pr01
```

```
1 ## Standard deviations (1, .., p=2):
2 ## [1] 3.535534 1.118034
3 ##
4 ## Rotation (n x k) = (2 x 2):
5 ##          PC1      PC2
6 ## [1,] -0.8944272 -0.4472136
7 ## [2,] -0.4472136  0.8944272
```

## 固有値

`Standard deviations` というのは、スコアの標準偏差  $\sqrt{\mu}$ .

```
1 pr01$sdev
```

```
1 ## [1] 3.535534 1.118034
```

二乗するとスコアの分散であり、 $X$ の分散・共分散行列  $M$ の固有値  $\mu$ .

```
1 pr01$sdev^2
```

```
1 ## [1] 12.50  1.25
```

この場合固有値は2個あるので、2つ出力。最初が第1主成分の、次のが第2主成分の固有値。第1主成分の方に分散が大きいことがわかる。

## 主成分のベクトル

`Rotation` は、主成分のベクトル  $(a_1, a_2)$ .

```
1 pr01$rotation
```

```
1 ##          PC1      PC2
2 ## [1,] -0.8944272 -0.4472136
3 ## [2,] -0.4472136  0.8944272
```

これも2種類あって、1列目が第1主成分で、2列目が第2主成分で、直交している。

第1主成分が  $(-2, -1)$ の方向で、第2主成分が  $(-1, 2)$ の方向で、ノルムが1となっている。

```
1 pr01$rotation[,1]^2
```

```
1 ## [1] 0.8 0.2
```

(横軸が  $X_1$ , 縦軸が  $X_2$  の場合) 左下方向と左上方向で、思っていたのと逆だが、どちら方向にスコアをとるかだけの問題で、本質的に違いはない。

## 各点のスコア

各点のスコアは, `pr01$x` で取り出せる。

```
1 pr01$x
```

```
1 ##          PC1      PC2
2 ## [1,]  0.000000  0.000000
3 ## [2,] -2.236068 -1.118034
4 ## [3,] -2.236068  1.118034
5 ## [4,] -4.472136 -1.118034
6 ## [5,] -4.472136  1.118034
7 ## [6,]  2.236068 -1.118034
8 ## [7,]  2.236068  1.118034
9 ## [8,]  4.472136 -1.118034
10 ## [9,]  4.472136  1.118034
```

スコアは  $a_1 x_1 + b_1 x_2$  のことだから, `x01` と `Rotation` を掛け合わせても得られる。

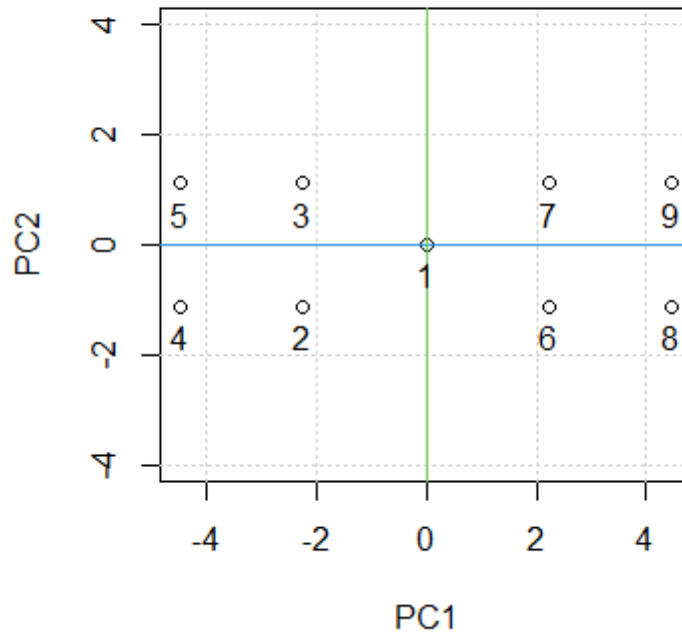
```
1 x01%*%pr01$rotation
```

```
1 ##          PC1      PC2
2 ## [1,]  0.000000  0.000000
3 ## [2,] -2.236068 -1.118034
4 ## [3,] -2.236068  1.118034
5 ## [4,] -4.472136 -1.118034
6 ## [5,] -4.472136  1.118034
7 ## [6,]  2.236068 -1.118034
8 ## [7,]  2.236068  1.118034
9 ## [8,]  4.472136 -1.118034
10 ## [9,]  4.472136  1.118034
```

## 各点のスコアのグラフ

第1主成分が横軸, 第2主成分が縦軸

```
1 plot(pr01$x, asp=1)
2 grid()
3 abline(h=0,col=4)
4 abline(v=0,col=3)
5 text(1:9,x=pr01$x[,1],y=pr01$x[,2],pos=1)
```



スコアのプロットをマイナス45度回して、さらに+-をひっくり返した形。  
 (青い線がもとの真ん中を通る直線)

例えば、2番目の点は、傾き1/2の直線（第1主成分方向）に垂線を下した点が(2,1)だから、第1主成分が $-\sqrt{2^2 + 1^2} = -\sqrt{5}$

```
1 | -sqrt(2^2+1^2)
```

```
1 | ## [1] -2.236068
```

第2主成分は、(2,1)の点から(0.5, -1)だけずらした分だから、 $-\sqrt{0.5^2 + (-1)^2} = -\sqrt{1.25}$

```
1 | -sqrt(0.5^2+(-1)^2)
```

```
1 | ## [1] -1.118034
```

スコアの分散は固有値

```
1 | var(pr01$x)%>%
2 | round(4)
```

```
1 | ##      PC1  PC2
2 | ## PC1 12.5  0.00
3 | ## PC2  0.0  1.25
```



# 主成分負荷量

個々の点の特徴を、2つの主成分に分けて見れることはわかった。

それでは $X_1$ という変数は、主成分から見てどんな変数か、 $X_2$ はどうか、というのが**主成分負荷量** (principal component loading)。

**因子負荷量** (factor loading)と呼ぶ人もいるが、因子分析の用語で気持ち悪い。

単に**負荷量** (loading) と呼ぶ場合も。

一般的な定義は

$$\left( \frac{a_1 \sqrt{\mu}}{\text{Var}(X_1)}, \frac{a_2 \sqrt{\mu}}{\text{Var}(X_2)} \right)$$

```
1 | vinv01<-1/apply(x01,2,sd)
2 | diag(vinv01)%*%pr01$rotation%*%diag(pr01$sdev)%>%
3 | round(4)
```

```
1 | ##      [,1]  [,2]
2 | ## [1,] -0.9877 -0.1562
3 | ## [2,] -0.8452  0.5345
```

これは $X$ とスコアとの相関係数。

$$\frac{\text{Cov}(X_1, a_1 X_1 + a_2 X_2)}{\sqrt{\text{Var}(X_1)} \sqrt{\text{Var}(a_1 X_1 + a_2 X_2)}} = \frac{a_1 \text{Var}(X_1) + a_2 \text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)} \sqrt{\mu}} = \frac{a_1 \mu}{\sqrt{\text{Var}(X_1)} \sqrt{\mu}}$$

3項目への変換は

$$\begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_1, X_2) & \text{Var}(X_2) \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \mu \begin{bmatrix} a_1 \\ a_2 \end{bmatrix}$$

からきている。

ふつうに相関係数を求めると。

```
1 | cor(x01,pr01$x)%>%
2 | round(4)
```

```
1 | ##      PC1      PC2
2 | ## [1,] -0.9877 -0.1562
3 | ## [2,] -0.8452  0.5345
```

定義した負荷量と一致する。

なので、主成分負荷量は $[-1, 1]$ の値をとる。

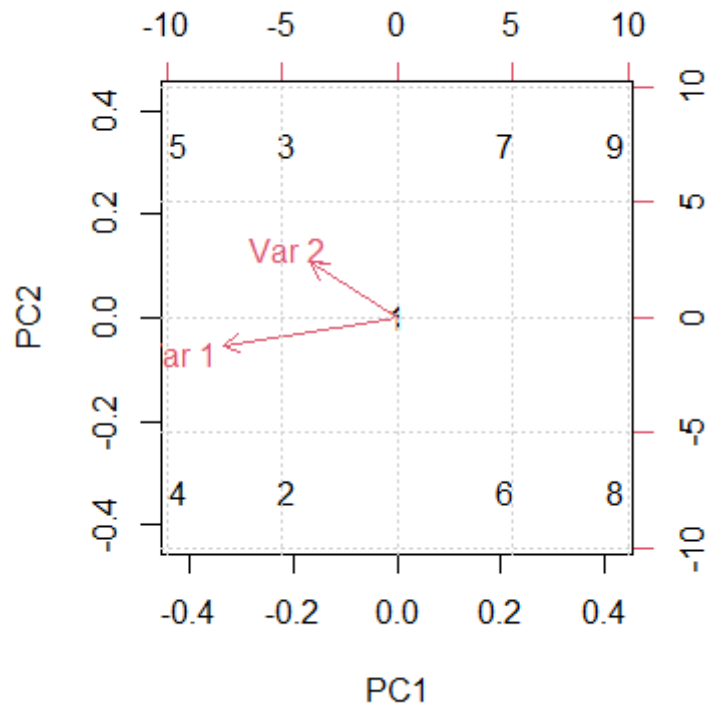
変数 $X_1$ は、第1主成分と-0.9877の相関があるが、第2主成分とは-0.1562しかない。

変数  $X_2$  は、第1主成分と-0.8452の相関があり、第2主成分とも0.5345の相関がある。

## biplot

主成分分析のスコアと主成分負荷量は、`biplot()` 関数で出力できる。

```
1 biplot(pr01,asp=1,scale=1)%>%
2 grid()
```



スコアと負荷量は全く別の数値だが、無理やり1つのグラフに入れてある。  
下と左の目盛がスコア用で、上と右が負荷量のための目盛。

1~9の数字が各点のスコアで、赤い矢印が負荷量・・・だと思ったが、何か数値が違う。

この図を理解するには、 $X$ の特異値分解をすればよくわかる。

## 特異値分解

$X$ の特異値分解  $X = UDV'$  は、Rを使えば `svd()` 関数で簡単にできる。

```
1 svd01<-svd(x01)
2 svd01%>%
3   lapply(.,function(x) round(x,4))
```

```
1 ## $d
2 ## [1] 10.0000 3.1623
3 ##
4 ## $u
5 ##      [,1] [,2]
6 ## [1,] 0.0000 0.0000
7 ## [2,] -0.2236 -0.3536
```

```

8 ## [3,] -0.2236 0.3536
9 ## [4,] -0.4472 -0.3536
10 ## [5,] -0.4472 0.3536
11 ## [6,] 0.2236 -0.3536
12 ## [7,] 0.2236 0.3536
13 ## [8,] 0.4472 -0.3536
14 ## [9,] 0.4472 0.3536
15 ##
16 ## $v
17 ##      [,1] [,2]
18 ## [1,] -0.8944 -0.4472
19 ## [2,] -0.4472 0.8944

```

```

1 u01<-svd01$u
2 d01<-svd01$d
3 v01<-svd01$v

```

$U$ も $V$ もノルム1の直交するベクトル

```
1 norm(u01,"2")
```

```
1 ## [1] 1
```

```
1 norm(v01,"2")
```

```
1 ## [1] 1
```

```

1 var(u01)%>%
2 round(4)

```

```

1 ##      [,1] [,2]
2 ## [1,] 0.125 0.000
3 ## [2,] 0.000 0.125

```

$DU$ で各点のスコアとなる。ただし、

$$D = \begin{bmatrix} d_1 & 0 \\ 0 & d_2 \end{bmatrix}$$

```

1 u01%*%diag(d01)%>%
2 round(4)

```

```

1 ##      [,1] [,2]
2 ## [1,] 0.0000 0.000
3 ## [2,] -2.2361 -1.118
4 ## [3,] -2.2361 1.118
5 ## [4,] -4.4721 -1.118
6 ## [5,] -4.4721 1.118
7 ## [6,] 2.2361 -1.118
8 ## [7,] 2.2361 1.118
9 ## [8,] 4.4721 -1.118
10 ## [9,] 4.4721 1.118

```

・・・であるはずなのだが、`biplot()` 関数では、スコアを以下で表し、

$$U \begin{bmatrix} d_1^{1-scale} & 0 \\ 0 & d_2^{scale} \end{bmatrix}$$

主成分負荷量を以下で表す。

$$V \begin{bmatrix} d_1^{scale} & 0 \\ 0 & d_2^{scale} \end{bmatrix}$$

ただし、*scale*は、 $[0, 1]$ の実数で、通常のスコアの定義は、 $scale = 0$ の場合である。

一方、`biplot()` のスコアは $scale = 1$ で出しているようだ。

この*scale*の値は、`biplot()` 関数の中で、引数 `scale=` で与えることができる。

## *scale* = 0の場合

スコアはUD

```
1 u01%*%diag(d01^(1-0))%>%
2 round(4)
```

```
1 ##      [,1] [,2]
2 ## [1,] 0.0000 0.000
3 ## [2,] -2.2361 -1.118
4 ## [3,] -2.2361  1.118
5 ## [4,] -4.4721 -1.118
6 ## [5,] -4.4721  1.118
7 ## [6,]  2.2361 -1.118
8 ## [7,]  2.2361  1.118
9 ## [8,]  4.4721 -1.118
10 ## [9,]  4.4721  1.118
```

・・・で、通常のスコアと一致する。

負荷量はV

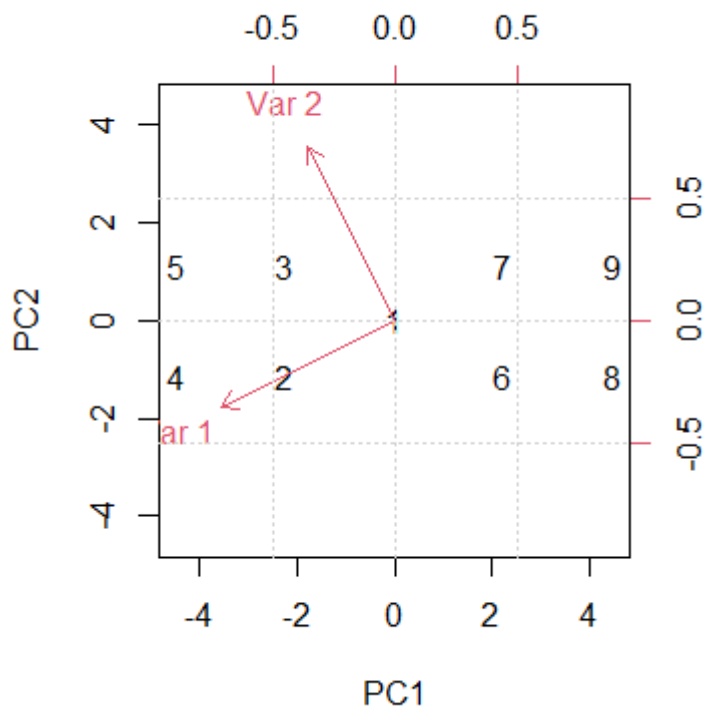
```
1 v01%*%diag(d01^0)%>%
2 round(4)
```

```
1 ##      [,1] [,2]
2 ## [1,] -0.8944 -0.4472
3 ## [2,] -0.4472  0.8944
```

で、これは主成分ベクトルそのもの。

`biplot` で確認

```
1 biplot(pr01, asp=1, scale=0)%>%
2 grid()
```



スコアの値は正確.

負荷量は, 定義した通りではあるのだが, 90度で直交して, 負荷量が相関係数を表していない. (だって, 主成分ベクトルだから...)

「 $X_1$ は, 第1主成分にくっついているので,  $X_1$ と第1主成分は相関が高い!」...とか, 負荷量を赤い矢印の傾きで判断したいのなら, 相関係数の傾きと同じになってほしい.

相関関係の傾きは

```
1 | c(-0.1562/-0.9877,0.5345/-0.8452)
```

```
1 | ## [1] 0.1581452 -0.6323947
```

scale = 0の傾きは

```
1 | c(-0.4472/-0.8944,0.8944/-0.4472)
```

```
1 | ## [1] 0.5 -2.0
```

...と, 全然違う. (だって主成分ベクトルだから)

## scale = 1の場合

スコアはU

```
1 | u01%*%diag(d01^(1-1))%>%
2 | round(4)
```

```

1 ##      [,1]  [,2]
2 ## [1,]  0.0000  0.0000
3 ## [2,] -0.2236 -0.3536
4 ## [3,] -0.2236  0.3536
5 ## [4,] -0.4472 -0.3536
6 ## [5,] -0.4472  0.3536
7 ## [6,]  0.2236 -0.3536
8 ## [7,]  0.2236  0.3536
9 ## [8,]  0.4472 -0.3536
10 ## [9,]  0.4472  0.3536

```

負荷量は,  $VD$

```

1 v01%*%diag(d01^1)%>%
2 round(4)

```

```

1 ##      [,1]  [,2]
2 ## [1,] -8.9443 -1.4142
3 ## [2,] -4.4721  2.8284

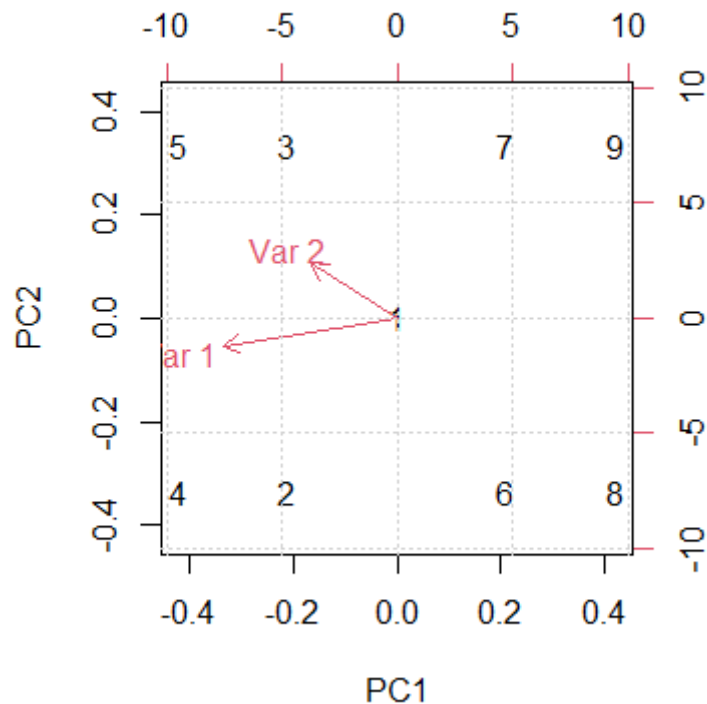
```

`biplot` で確認すると

```

1 biplot(pr01,asp=1,scale=1)%>%
2 grid()

```



負荷量の傾きは・・・

```

1 c(-1.4142/-8.9443,2.8284/-4.4721)

```

```

1 ## [1]  0.1581119 -0.6324546

```

相関係数の傾きとほぼ一緒。（値は違うけど）

$$scale = 0.5$$

$D$ の各要素の平方根をとったものを $D^{\frac{1}{2}}$ と表すと・・・

スコアは $UD^{\frac{1}{2}}$

```
1 u01%*%diag(d01^(1-0.5))%>%
2 round(4)
```

```
1 ##      [,1]  [,2]
2 ## [1,]  0.0000  0.0000
3 ## [2,] -0.7071 -0.6287
4 ## [3,] -0.7071  0.6287
5 ## [4,] -1.4142 -0.6287
6 ## [5,] -1.4142  0.6287
7 ## [6,]  0.7071 -0.6287
8 ## [7,]  0.7071  0.6287
9 ## [8,]  1.4142 -0.6287
10 ## [9,]  1.4142  0.6287
```

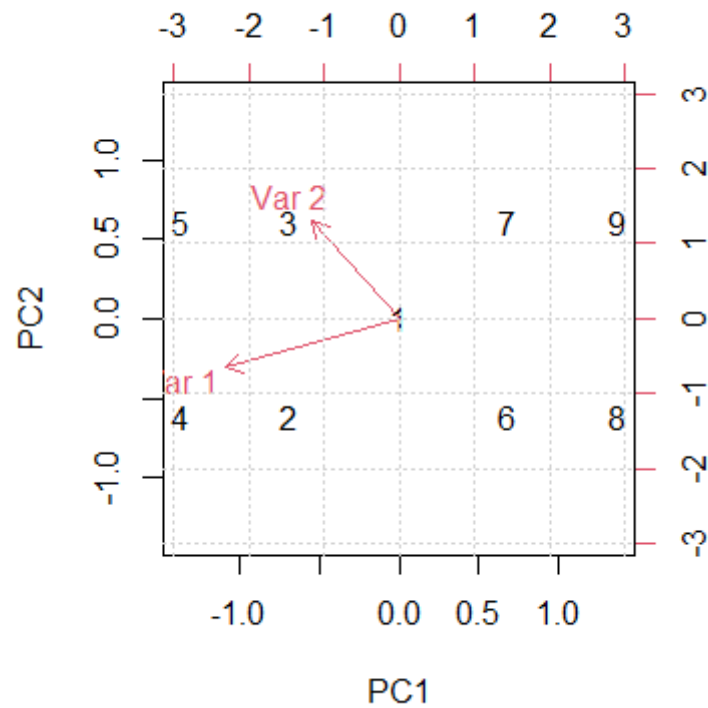
負荷量は $VD^{\frac{1}{2}}$

```
1 v01%*%diag(d01^0.5)%>%
2 round(4)
```

```
1 ##      [,1]  [,2]
2 ## [1,] -2.8284 -0.7953
3 ## [2,] -1.4142  1.5905
```

確認する

```
1 biplot(pr01,asp=1,scale=0.5)%>%
2 grid()
```



負荷量の傾きは

```
1 | c(-0.7953/-2.8284,1.5905/-1.4142)
```

```
1 | ## [1] 0.2811837 -1.1246641
```

・・・微妙.

結局, 負荷量を相関係数として理解するなら,  $scale = 1$ でbiplotを描くのが最もそれに近い. ただし, 目盛は違うので要注意・・・ということ.

## 寄与率

スコアの分散が固有値だが, 分散の和における個々の分散 (固有値) の割合は**寄与率**と呼ばれる.

固有値が第2主成分まで, 2つある場合 ( $\mu_1, \mu_2$ )

第1主成分の寄与率は

$$\frac{\mu_1}{\mu_1 + \mu_2}$$

2変数しかない場合,  $\mu_1 + \mu_2$ が $X$ の分散 ( $Var(X_1) + Var(X_2)$ ) と一致するので, この寄与率は第1主成分が $X$ の分散を説明する割合ということになる.

スコアの分散は

```
1 | var(pr01$x)%>%
2 | round(4)
```



```
1 ##      PC1  PC2
2 ## PC1 12.5 0.00
3 ## PC2  0.0 1.25
```

固有値は

```
1 pr01$sdev^2
```

```
1 ## [1] 12.50  1.25
```

寄与率は

```
1 pr01$sdev^2/sum(pr01$sdev^2)%>%
2 round(4)
```

```
1 ## [1] 0.90909091 0.09090909
```

第1主成分の寄与率が90.9%で、第2主成分の寄与率が9.09%ということ。

寄与率は、`summary()` 関数で出力できる。

```
1 summary(pr01)
```

```
1 ## Importance of components:
2 ##                PC1    PC2
3 ## Standard deviation  3.5355 1.11803
4 ## Proportion of Variance 0.9091 0.09091
5 ## Cumulative Proportion 0.9091 1.00000
```

`Proportion of Variance` が寄与率。それを大きい方から順に足していったのが累積寄与率。`Cumulative Proportion` がそれ。この場合2変数なので、主成分は高々2個。なので、第2主成分で累積寄与率は1になる。

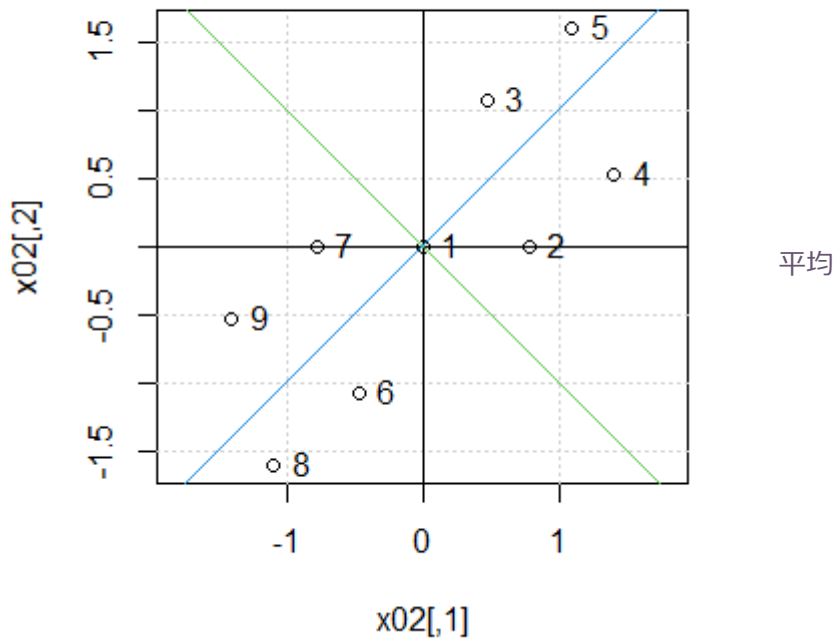
## 標準化

主成分分析は、 $X$ の各変数の単位をどうとるかで、出力される値が変わってしまうので、（よほど変動を変えたくない場合を除き）通常は**全ての変数の分散を1**に標準化してから主成分分析を行うことが一般的なようだ。

## $X$ を標準化する。

`scale()`

```
1 x02<-scale(x01)
2 plot(x02,asp=1)
3 grid()
4 abline(h=0)
5 abline(v=0)
6 abline(a=0,b=1,col=4)
7 abline(a=0,b=-1,col=3)
8 text(1:9,x=x02[,1],y=x02[,2],pos=4)
```



```
1 apply(x02,2,mean)
```

```
1 ## [1] 0 0
```

分散

```
1 var(x02)
```

```
1 ##          [,1]      [,2]
2 ## [1,] 1.0000000 0.7513055
3 ## [2,] 0.7513055 1.0000000
```

## 主成分分析

```
1 pr02<-prcomp(x02)
2 pr02
```

```
1 ## Standard deviations (1, .., p=2):
2 ## [1] 1.3233690 0.4986928
3 ##
4 ## Rotation (n x k) = (2 x 2):
5 ##          PC1      PC2
6 ## [1,] -0.7071068 -0.7071068
7 ## [2,] -0.7071068  0.7071068
```

2変数しかない場合、主成分のベクトルは、2本の45度線。

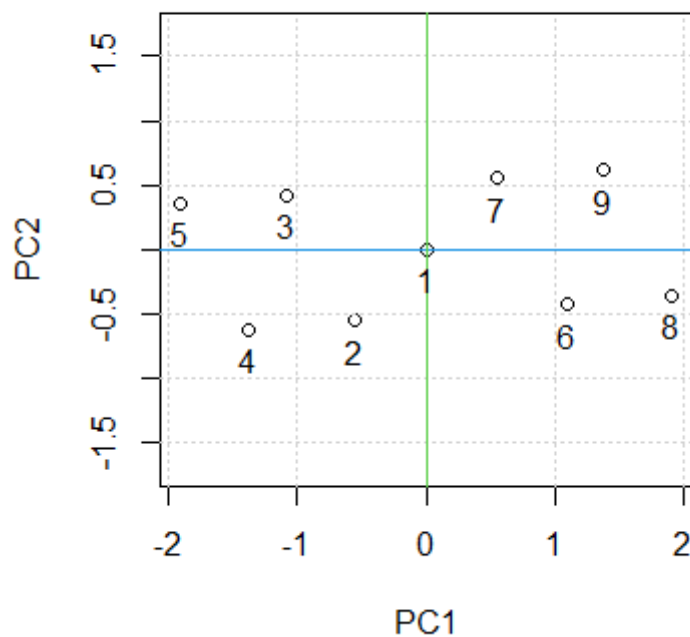
事前にXを標準化しなくても `prcomp()` 関数に引数 `scale=T` を与えても同じ結果となる。

```
1 prcomp(x01,scale=T)
```

```
1 ## Standard deviations (1, ..., p=2):  
2 ## [1] 1.3233690 0.4986928  
3 ##  
4 ## Rotation (n x k) = (2 x 2):  
5 ##           PC1      PC2  
6 ## [1,] -0.7071068 -0.7071068  
7 ## [2,] -0.7071068  0.7071068
```

## スコアのプロット

```
1 plot(pr02$x,asp=1)  
2 grid()  
3 abline(h=0,col=4)  
4 abline(v=0,col=3)  
5 text(1:9,x=pr02$x[,1],y=pr02$x[,2],pos=1)
```



## 主成分負荷量

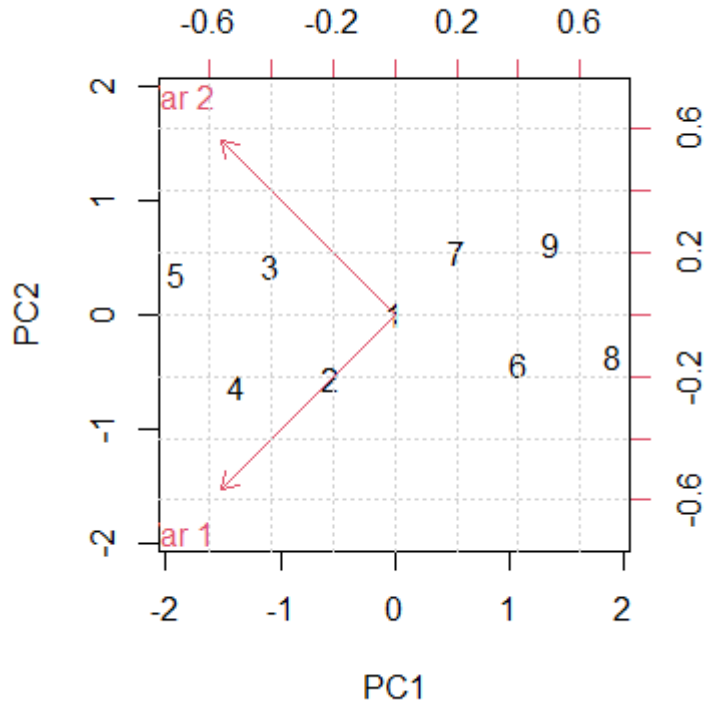
```
1 cor(x02,pr02$x)
```

```
1 ##           PC1      PC2  
2 ## [1,] -0.9357632 -0.3526291  
3 ## [2,] -0.9357632  0.3526291
```

# biplot

*scale = 0*の場合

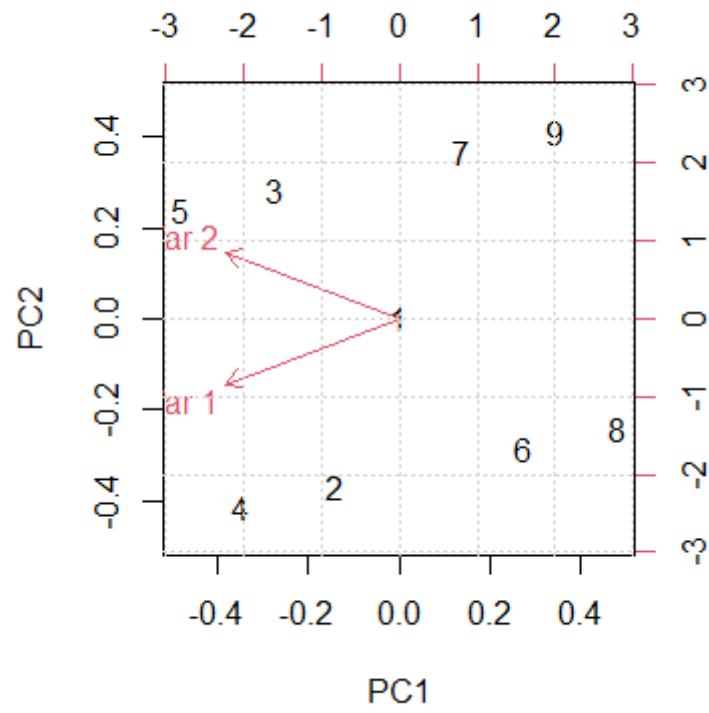
```
1 biplot(pr02,scale = 0)
2 grid()
```



負荷量（赤線）の値が、なんかちょっと違う気がするが・・・この辺りは謎.

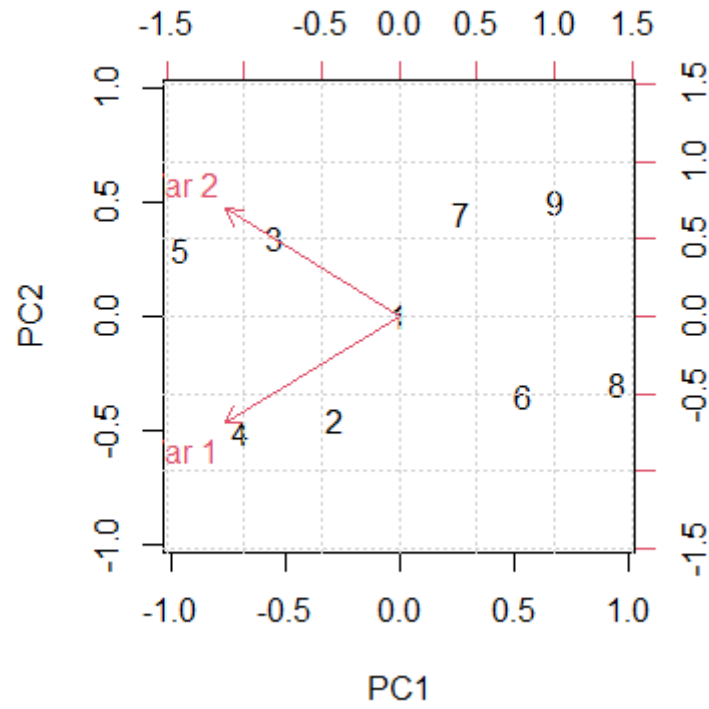
*scale = 1*の場合

```
1 biplot(pr02,scale = 1)
2 grid()
```



*scale = 0.5*

```
1 biplot(pr02,scale = 0.5)
2 grid()
```



# 寄与率

```
1 summary(pr02)
```

```
1 ## Importance of components:
2 ##                PC1    PC2
3 ## Standard deviation  1.3234 0.4987
4 ## Proportion of Variance 0.8757 0.1244
5 ## Cumulative Proportion 0.8757 1.0000
```

第1主成分の寄与率が、標準化していない場合よりも小さくなった。

主成分分析の結果は変数のスケールに応じて変わるということ。つまり同じデータでも、例えばキログラムで測るのか、それともトンで測るのかでも結果が変わるということ。

それはまずい、ということで、主成分分析では各変数の分散を1に揃えるのが一般的。

ただし、合計に意味があるような場合。例えば、入試の成績などでは、各科目の分散を1に揃えない方が、分析の意味づけがやりやすくなるはず、たぶん。

---

1. 有馬哲・石村貞夫『多変量解析のはなし』東京図書,1987 →