

Partial Least Squares (PLS; 部分最小二乗法, 偏最小二乗法)

PLSってどのくらい使えるのか？

1. PLSは $(Xu)'Yv$ を最大化するように u^* と v^* を決める。
2. この u^* と v^* は、 $X'Y$ の特異値分解から得られる。
3. 推定の考え方がOLSとは違うので、 X と Y にかく乱項が加わった場合、係数の不偏性は怪しい。
4. PLSの実際の利用は、多数の変数を持つ X から Y の説明力の高い X の主成分を抽出し、それと相関の高い X を特定すること。
5. 変数間の係数に関心がある場合は、PLSで特定した X でPCRをするのがよさそう。
6. 実際のPLSは、 $X'Y$ の特異値分解という単純なものではなく、主成分を1つずつ算出していく複雑なアルゴリズムが複数ある。
7. Rのパッケージplsにある `pls()` 関数で出力される Y のスコア (`Yscores`) は X のスコアに合わせてスケールを調整しているので要注意。

説明変数 X で Y に回帰させる一般的な方法がOLSなのだが、 X に相関があると回帰がうまくいかない。

そこで、 X を主成分分析して、相関する変数をまとめた主成分を抽出する主成分回帰分析 (PCR) というのがある。しかし、この X の主成分分析には Y の情報を含まないので、必ずしも主成分分析の結果が意味のあるものでないかもしれない。

そこでそこで、 Y の情報も含めて主成分分析的な回帰をやろうというのがPLS。

もともと心理学の分野で開発されたらしいが、最近は化学の分野で多用されているらしい。特定の反応を説明するのに、サンプル数が数十個しかないのに対して、説明変数の候補が100個以上あったりする場合に、説明変数の候補を探し出すのに使われているようだ。

我々が行う購買行動調査のようなアンケートも、性別や年齢などの属性で説明できるような多純なものではなく、様々な消費性向のようなことを重ねて聞くから、どれとどれがどれと関連しているかを突き止めるのが大変。

PLSはあくまで連続変数間の関係なので、アンケート分析にどれだけ使えるかわからないが、なんかの参考にはなるだろう。PLSは被説明変数に複数の変数をとれるというのも魅力的だ。

気になったので勉強しようとしてみた。ところが、どうも教科書が無いのか？ネット上の解説をあさってみた。上記のような話が出てきたが、具体的に何をしているかがよくわからない。たくさん抽出できる主成分の中からどうやって主成分を抽出するかという手順が説明してある。しかし、その前に**そもそもどんな回帰をしているのか？**その説明が飛んでいる。PLSの場合、この手順が醍醐味なのはよくわかったが、言ってしまうと、そこは計算機の仕事でしょ、ということになる。それよりも、結果として出力される値の意味が知りたい。少なくとも分かった気になりたい。


```
14 | x1<-rbind(x1_0,x1_1)
```

平均は2変数とも0

```
1 | apply(x1_1,2,mean)%>%  
2 | round(4)
```

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

分散・共分散は

```
1 | var(x1)
```

```
1 | ##      [,1] [,2]  
2 | ## [1,]  1.0  0.8  
3 | ## [2,]  0.8  1.0
```

・・・で、 X のデータは、うまいことできた！

プロット図に書いてみよう

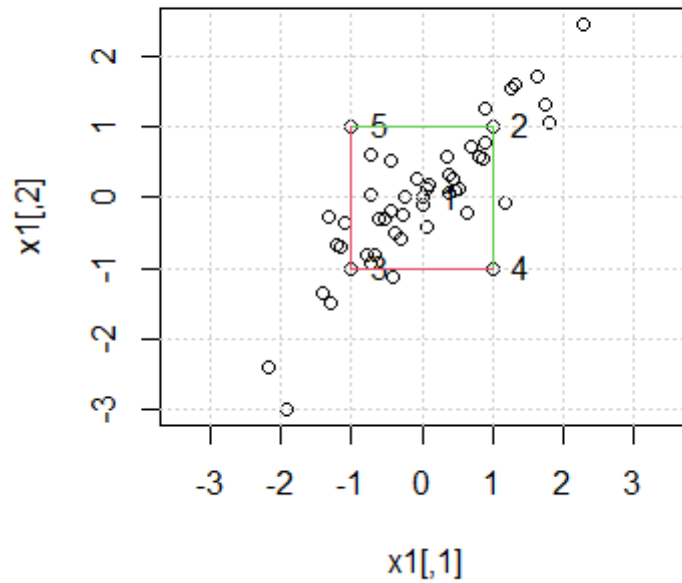
データの変換をプロット図で追いかけられるように、2~5行目のデータは、特殊な点を無理やり入れてある。そのをつなぐ線を描きたい。

その関数を定義

```
1 | f_grid1<-function(x1){  
2 |   segments(x1[2,1],x1[2,2],x1[4,1],x1[4,2],col=3)  
3 |   segments(x1[2,1],x1[2,2],x1[5,1],x1[5,2],col=3)  
4 |   segments(x1[3,1],x1[3,2],x1[4,1],x1[4,2],col=2)  
5 |   segments(x1[3,1],x1[3,2],x1[5,1],x1[5,2],col=2)  
6 | }
```

X のプロット図。1が原点、2が(1,1)、3が(-1,-1)、4が(1,-1)、5が(-1,1)。

```
1 | plot(x1,asp=1)  
2 | grid()  
3 | text(x=x1[1:5,1],y=x1[1:5,2],1:5,pos = 4)  
4 | f_grid1(x1[1:5,])
```



Xの主成分のスコアの計算

Xの主成分分析をすると・・・

```
1 prc_x1<-prcomp(x1)
```

固有ベクトルは、第1主成分が右上がりの45度線、第2主成分がその直交（右下がり45度線）

```
1 prc_x1$rotation
```

```
1 ##          PC1          PC2
2 ## [1,] 0.7071068  0.7071068
3 ## [2,] 0.7071068 -0.7071068
```

固有値は

```
1 prc_x1$sdev^2>%
2 round(4)
```

```
1 ## [1] 1.8 0.2
```

スコアを求める。

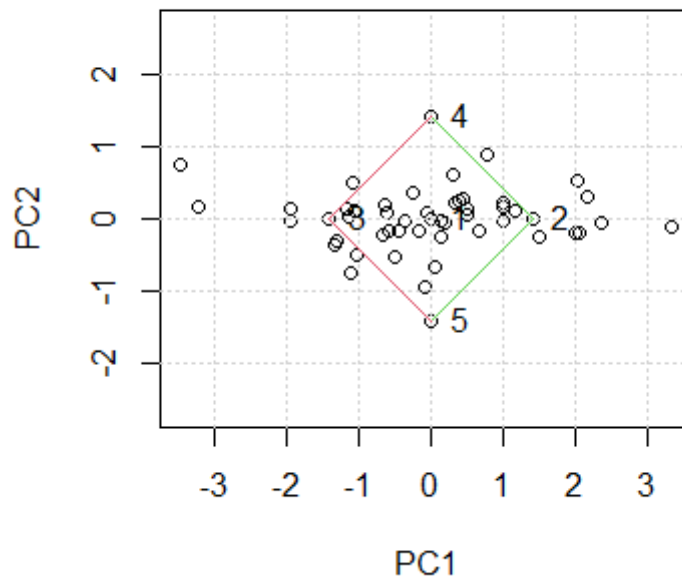
```
1 Px1<-prc_x1$x
```

スコアのプロット図：横軸が第1主成分、縦軸が第2主成分。
Xのプロットを45度度マイナスにした形となる。

```

1 plot(Px1,asp=1)
2 grid()
3 text(x=Px1[1:5,1],y=Px1[1:5,2],1:5,pos=4)
4 f_grid1(Px1)

```



X のプロット図の45度線でスコアを測っていることがわかる。

X の2変数とも標準化してしまうと、どうあがいても、固有ベクトルの方向は、45度線になってしまう。

第1主成分が (x_1, x_2) の方のプロット図の右上がり45度線。

第2主成分が同図の右下がり45度線

つまり、固有ベクトルは、45度($\pi/4$)線で、長さ(ノルム)1なので・・・

$$\begin{pmatrix} \cos \frac{\pi}{4}, \cos \frac{-\pi}{4} \\ \sin \frac{\pi}{4}, \sin \frac{-\pi}{4} \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}} \end{pmatrix}$$

ということ。

各主成分で符号が逆になる場合があるが、±どっち向きにスコアを測るかが違うだけで、本質的に同じ。

なので、 $1/\sqrt{2} = 0.7071... = a$ と置いておくと・・・

```

1 a<-1/sqrt(2)

```

スコアは、 $ax_1 + ax_2$ (第1主成分) , $ax_1 - ax_2$ (第2主成分)

```

1 cbind(a*x1[,1]+a*x1[,2],a*x1[,1]-a*x1[,2])%>%
2   round(4)%>%
3   head()

```

```

1 ##      [,1]  [,2]
2 ## [1,] 0.0000 0.0000
3 ## [2,] 1.4142 0.0000
4 ## [3,] -1.4142 0.0000
5 ## [4,] 0.0000 1.4142
6 ## [5,] 0.0000 -1.4142
7 ## [6,] -1.0460 0.1174

```

ちゃんと主成分分析のスコアと一致しますね。

```

1 Px1%>%
2   round(4)%>%
3   head()

```

```

1 ##      PC1    PC2
2 ## [1,] 0.0000 0.0000
3 ## [2,] 1.4142 0.0000
4 ## [3,] -1.4142 0.0000
5 ## [4,] 0.0000 1.4142
6 ## [5,] 0.0000 -1.4142
7 ## [6,] -1.0460 0.1174

```

固有値は、それぞれ**スコアの分散**なんですね、これが！

さらに、 $Var(x_1) = Var(x_2) = 1$ なので、

第1主成分の固有値は、この場合、 $1 + Cov(x_1, x_2) = 1.8$

第2主成分の固有値は、この場合、 $1 - Cov(x_1, x_2) = 0.2$

```

1 var(prc_x1$x)%>%
2   round(4)

```

```

1 ##      PC1 PC2
2 ## PC1 1.8 0.0
3 ## PC2 0.0 0.2

```

Yを作成する

方針

XとYの直接的な規定関係ではなく、複数のXのスコアと複数のYのスコアとの関係があるものと考えよう。

観測できる変数（顕在変数）XとYは、観測できない変数（潜在変数、XのスコアとYのスコア）がノイズをもって観測されたものとする。

最も簡単な例として、Xが上記2変数で、Yも2変数。

(Yが1変数という場合もあるが、PLSが2変数以上を被説明変数にできるというのだから、そこは2変数としよう。)

X と Y は、それぞれ2変数とも平均0、標準偏差1に標準化.

ひとまず、 x と Y にノイズは一切入らないとする.

Y のスコアを P_{y1}, P_{y2} , X のスコアを P_{x1}, P_{x2} と表し...

$$P_{y1} = 0.7P_{x1} + 1.3P_{x2}$$

$$P_{y2} = 0.346489P_{x1} - 1.679139P_{x2}$$

...で Y のスコアが決まるとする.

Y は、スコアから45度逆回転させればよいのだから.

$$P_{y1} = \frac{1}{\sqrt{2}}y_1 + \frac{1}{\sqrt{2}}y_2$$

$$P_{y2} = \frac{1}{\sqrt{2}}y_1 - \frac{1}{\sqrt{2}}y_2$$

で、 $Y = (y_1, y_2)$ を求めることができる.

手順

まず y_1 の2つの主成分のスコア P_{y1}, P_{y2} を求める.

第1主成分 P_{y1} が、 X の主成分のスコアと次のような関係にあるとすると... (0.7, 1.3という値は適当に決めた)

$$P_{y1} = 0.7P_{x1} + 1.3P_{x2}$$

```
1 Py1_1<-0.7*Px1[,1]+1.3*Px1[,2]
```

この時、 P_{y1} の分散は

```
1 var(Py1_1)
```

```
1 ## [1] 1.22
```

Y の第2主成分

$$P_{y2} = b_1P_{x1} + b_2P_{x2}$$

は、これと直交するから

$$P'_{y1}P_{y2} = 0$$

となるように P_{y2} を決めたい.

$$(0.7P_{x1} + 1.3P_{x2})'(b_1P_{x1} + b_2P_{x2}) = 0$$

$$(0.7 \quad 1.3) \begin{pmatrix} P'_{x1} \\ P'_{x2} \end{pmatrix} (P_{x1} \quad P_{x2}) \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = 0$$

$$(0.7 \quad 1.3) \begin{pmatrix} 90 & 0 \\ 0 & 10 \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} = 0$$

$$0.7 \times 90b_1 + 1.3 \times 10b_2 = 0$$

さらに、2つのスコアの分散の合計を2にしたいから

$$\text{Var}(P_{y2}) = 2 - 1.22 = 0.78$$

$$1.8b_1^2 + 0.2b_2^2 = 0.78$$

b_2/b_1 は

```
1 b2pb1<- -0.7*90/(1.3*10)
```

b_1 は

```
1 b1<-sqrt(0.78/(1.8+0.2*b2pb1^2))
```

b_2 は

```
1 b2<-b2pb1*b1
```

つまり、 (b_1, b_2) は

```
1 c(b1,b2)
```

```
1 ## [1] 0.346489 -1.679139
```

P_{y1} の係数と合わせて `b0` に代入しておこう。

```
1 b0<-matrix(c(0.7,1.3,b1,b2),nrow=2)
```

$P_{y2} = b_1 P_{x1} + b_2 P_{x2}$ が完成!

```
1 Py1_2<-b1*Px1[,1]+b2*Px1[,2]
```

P_{y1} と合わせて、 $P_y = (P_{y1}, P_{y2})$.

```
1 Py1<-cbind(Py1_1,Py1_2)
```

直交しているか確認。

```
1 var(Py1)%>%  
2 round(4)
```

```
1 ##      Py1_1 Py1_2  
2 ## Py1_1  1.22  0.00  
3 ## Py1_2  0.00  0.78
```


ちゃんと直交しているし、分散も、合わせて2になっている。
 P_y ができた。

次に、 Y を求める。

これは簡単。 P_y を45度逆回転させればよいだけ。

```
1 a1<-matrix(c(a,a,a,-a),ncol=2)
2 y1<-Py1%*%solve(a1)
```

Y の平均は、

```
1 apply(y1,2,mean)%>%
2 round(4)
```

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

Y の分散・共分散

```
1 var(y1)%>%
2 round(4)
```

```
1 ##      [,1] [,2]
2 ## [1,] 1.00 0.22
3 ## [2,] 0.22 1.00
```

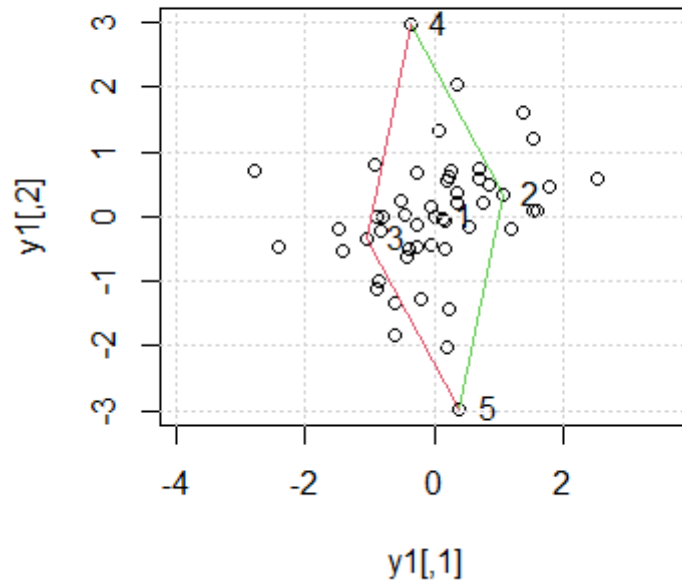
ちゃんと標準化されてる！

Y のデータ完成。

Y のプロット

Y の分布を確認する

```
1 plot(y1,asp=1)
2 grid()
3 text(x=y1[1:5,1],y=y1[1:5,2],1:5,pos = 4)
4 f_grid1(y1)
```

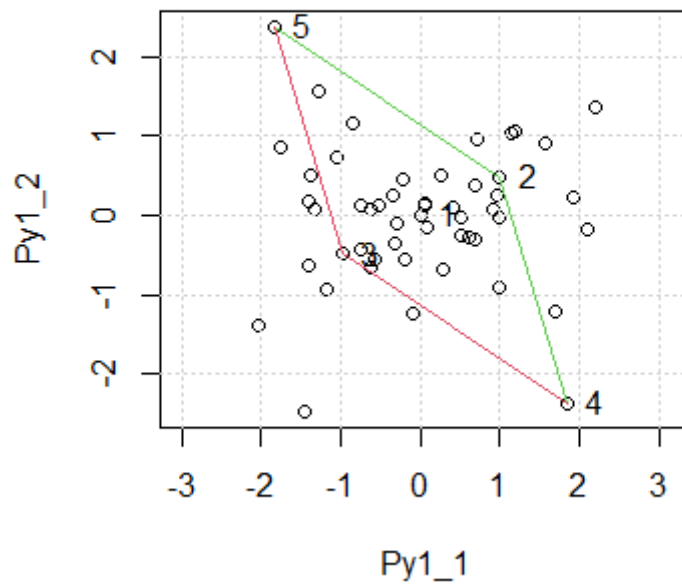


Y のスコアのプロット. Y のプロットを45度マイナスにした形となる.

```

1 plot(Py1,asp=1)
2 grid()
3 text(x=Py1[1:5,1],y=Py1[1:5,2],1:5,pos=4)
4 f_grid1(Py1)

```



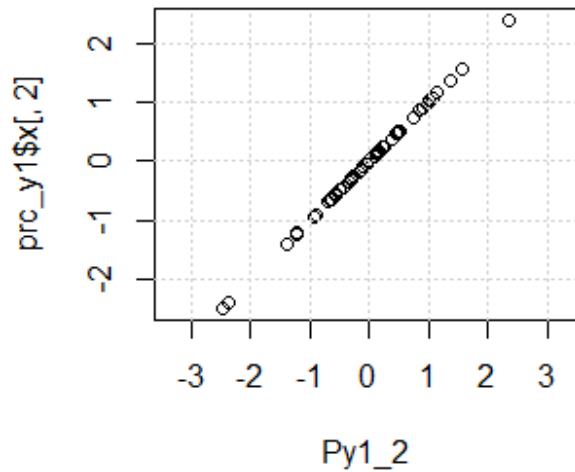
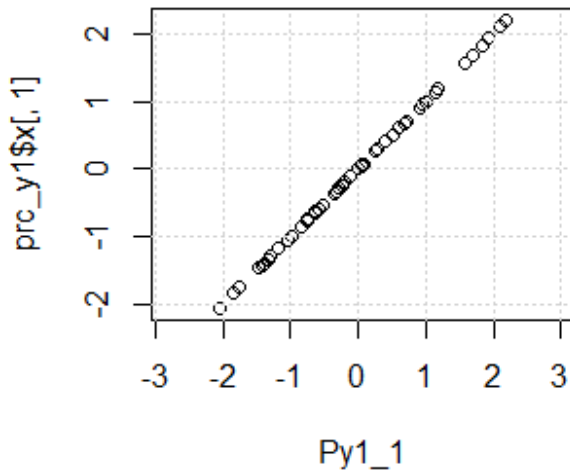
・・・という感じ.

(P_{y1}, P_{y2}) も Y の主成分となっていることを確認しておく.

```

1 prc_y1<-prcomp(y1)
2 par(mfrow=c(1,2))
3 plot(prc_y1$x[,1]~Py1_1,asp=1)
4 grid()
5 plot(prc_y1$x[,2]~Py1_2,asp=1)
6 grid()

```



・・・と、まあ、データの作成に手間取ったが、通常は、「こんなデータがありました！」ということで、ここからが分析のスタートとなる。

OLS回帰とPCR

関係式

X と Y の関係式は以下

$X \rightarrow$ 主成分 P_x

$$P_{x1} = \frac{1}{\sqrt{2}}x_1 + \frac{1}{\sqrt{2}}x_2$$

$$P_{x2} = \frac{1}{\sqrt{2}}x_1 - \frac{1}{\sqrt{2}}x_2$$

$$A = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$

とすると、

$$P_x = XA$$

(P_x は51行×2列, X も51行2列, A は2行×2列)

$Y \rightarrow$ 主成分 P_y

$$P_{y1} = \frac{1}{\sqrt{2}}y_1 + \frac{1}{\sqrt{2}}y_2$$

$$P_{y2} = \frac{1}{\sqrt{2}}y_1 - \frac{1}{\sqrt{2}}y_2$$

これも・・・

$$P_y = YA$$

(P_y は51行×2列, y も51行2列, A は2行×2列)

主成分 $P_x \rightarrow$ 主成分 P_y

$$P_{y1} = 0.7P_{x1} + 1.3P_{x2}$$

$$P_{y2} = 0.346489P_{x1} - 1.679139P_{x2}$$

$$B = \begin{pmatrix} 0.7 & 1.3 \\ 0.346489 & -1.679139 \end{pmatrix}$$

とすると,

$$P_y = P_x B$$

(P_y は51行×2列, P_x も51行2列, B は2行×2列)

OLS回帰

作成した X と Y でOLS回帰してみる.

y_1 への回帰は・・・

```
1 lm(y1[,1]~x1-1)%>%
2 summary()
```

```
1 ## Warning in summary.lm(.): essentially perfect fit: summary may be unreliable
2
3 ##
4 ## Call:
5 ## lm(formula = y1[, 1] ~ x1 - 1)
6 ##
7 ## Residuals:
8 ##      Min       1Q   Median       3Q      Max
9 ## -4.951e-16 -6.125e-17  2.570e-18  9.978e-17  1.918e-15
10 ##
11 ## Coefficients:
12 ##      Estimate Std. Error  t value Pr(>|t|)
13 ## x11  3.337e-01  7.162e-17  4.659e+15 <2e-16 ***
14 ## x12  7.128e-01  7.162e-17  9.953e+15 <2e-16 ***
15 ## ---
16 ## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
17 ##
18 ## Residual standard error: 3.039e-16 on 49 degrees of freedom
```

```
19 ## Multiple R-squared: 1, Adjusted R-squared: 1
20 ## F-statistic: 2.708e+32 on 2 and 49 DF, p-value: < 2.2e-16
```

データにかく乱項を全く入れなかったので推定誤差は0。以下同様なので、係数だけ見る。

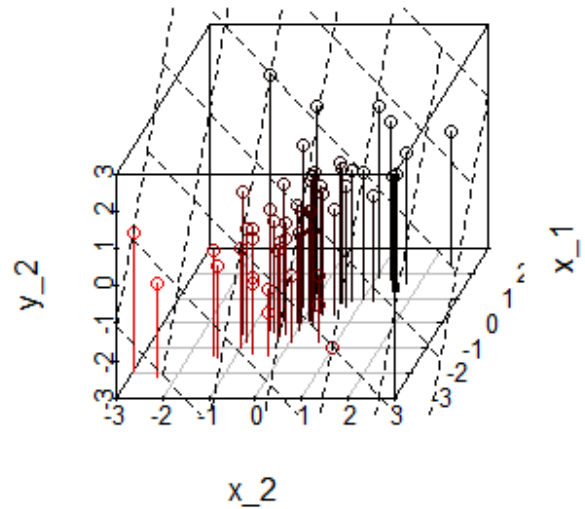
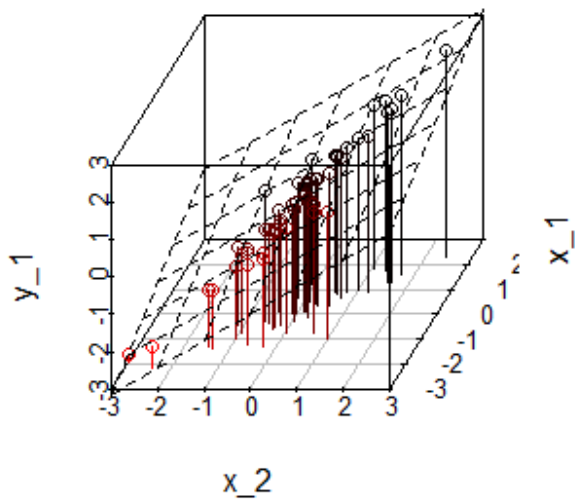
```
1 lm(y1~x1-1)%>%
2   coef()%>%
3   round(4)
```

```
1 ##      [,1]  [,2]
2 ## x11 0.3337 1.6663
3 ## x12 0.7128 -1.3128
```

(被説明変数を2変数にすると、それぞれの回帰をいっぺんにやってくれる。1列目が $y_1[,1]$ に対する回帰の結果、2列目が $y_1[,2]$ に対する回帰の結果)

一見、 y_1 に対して x_1 が1.38、 x_2 がその半分ぐらいの強さで影響しているように見えるが、プロットしてみると・・・

```
1 d1_1_OLS<-data.frame(x_2=x1[,2],
2                     x_1=x1[,1],
3                     y_1=y1[,1])
4
5 d1_2_OLS<-data.frame(x_2=x1[,2],
6                     x_1=x1[,1],
7                     y_2=y1[,2])
8
9 par(mfrow=c(1,2))
10
11 g3d1_1_OLS<-scatterplot3d(d1_1_OLS,
12                          type="h",
13                          angle=60,
14                          highlight.3d = T)
15 fit1_1_OLS<-lm(y_1~x_2+x_1-1,d1_1_OLS)
16 g3d1_1_OLS$plane3d(fit1_1_OLS)
17
18 g3d1_2_OLS<-scatterplot3d(d1_2_OLS,
19                          type="h",
20                          angle=60,
21                          highlight.3d = T)
22 fit1_2_OLS<-lm(y_2~x_2+x_1-1,d1_2_OLS)
23 g3d1_2_OLS$plane3d(fit1_2_OLS)
```



X の相関が強く、平面を回帰するのは無理な感じがするが、誤差はないので、これはこれ。

$$P_x = XA$$

$$P_y = YA$$

$$P_y = P_x B$$

という関係にあるので、 Y に X を直接回帰したら、

$$Y = P_y A^{-1} = P_x B A^{-1} = X A B A^{-1}$$

であり、 X の係数は $A B A^{-1}$ ということだ。

これを上の例で計算してみると・・・

```
1 a1%*%b0%*%solve(a1)%>%
2 round(4)
```

```
1 ##      [,1]  [,2]
2 ## [1,] 0.3337  1.6663
3 ## [2,] 0.7128 -1.3128
```

OLSの結果と一致する。

主成分回帰 (PCR)

X のスコアと Y との回帰で、主成分回帰にあたる。

```
1 lm(y1~Px1-1)%>%
2 coef()
```

```

1 ##           [,1]    [,2]
2 ## Px1PC1  0.7399795 0.249970
3 ## Px1PC2 -0.2680918 2.106569

```

$Y = P_y A^{-1} = P_x B A^{-1}$ であり, P_x の係数は BA^{-1} .

上の例だと・・・

```

1 b0%*%solve(a1)

```

```

1 ##           [,1]    [,2]
2 ## [1,]  0.7399795 0.249970
3 ## [2,] -0.2680918 2.106569

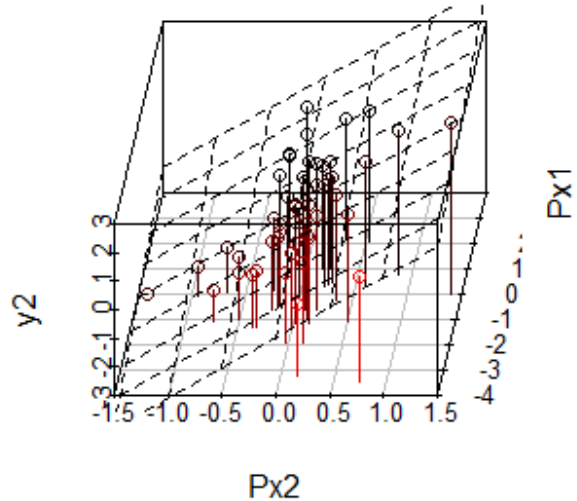
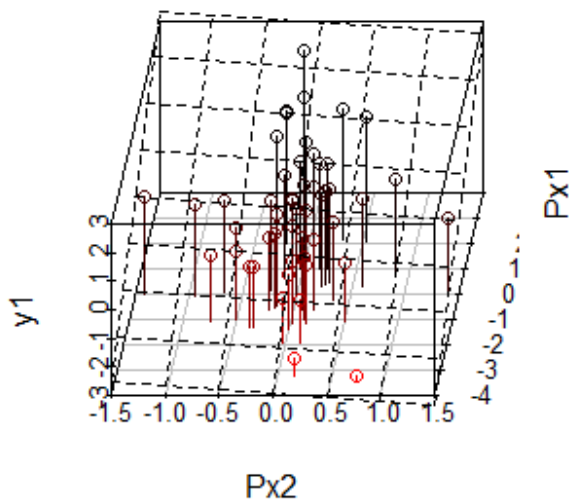
```

一致する.

```

1 d1_1_PCR<-data.frame(Px2=Px1[,2],
2                     Px1=Px1[,1],
3                     y1=y1[,1])
4 d1_2_PCR<-data.frame(Px2=Px1[,2],
5                     Px1=Px1[,1],
6                     y2=y1[,2])
7
8 par(mfrow=c(1,2))
9
10 g3d1_1_PCR<-scatterplot3d(d1_1_PCR,
11                          type="h",
12                          angle=80,
13                          highlight.3d = T)
14 fit1_1_PCR<-lm(y1~Px2+Px1-1,d1_1_PCR)
15 g3d1_1_PCR$plane3d(fit1_1_PCR)
16
17 g3d1_2_PCR<-scatterplot3d(d1_2_PCR,
18                          type="h",
19                          angle=80,
20                          highlight.3d = T)
21 fit1_2_PCR<-lm(y2~Px2+Px1-1,d1_2_PCR)
22 g3d1_2_PCR$plane3d(fit1_2_PCR)

```



この場合、 y_1 は' Px_1 '方向の傾きだけで、 y_2 は' Px_2 '方向の傾きだけで説明できるような形となっている。

PLS的な回帰

主成分分析のスコアどうしでOLS回帰してみる。

```
1 lm(Py1~Px1-1)%>%
2   coef()%>%
3   round(4)
```

```
1 ##      Py1_1  Py1_2
2 ## Px1PC1  0.7  0.3465
3 ## Px1PC2  1.3 -1.6791
```

$P_y = P_x B$ であり、 X の係数は定義通りで、 B 。上の例だと・・・

```
1 b0
```

```
1 ##      [,1]      [,2]
2 ## [1,]  0.7  0.346489
3 ## [2,]  1.3 -1.679139
```

Y の第1主成分のスコアを X の第1主成分のスコアのみで説明していて、 Y の第2主成分のスコアを X の第2主成分のスコアのみで説明している。

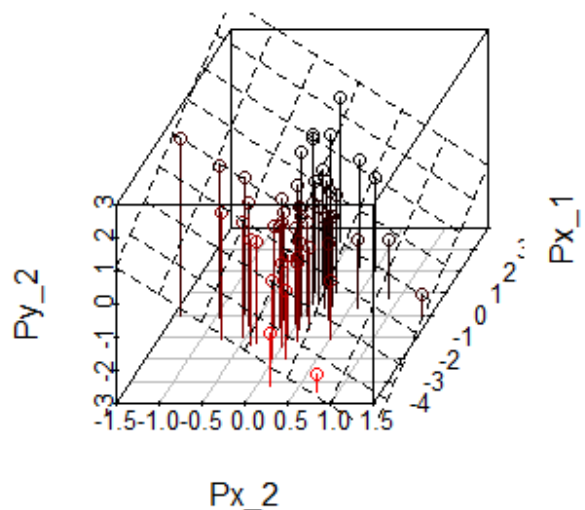
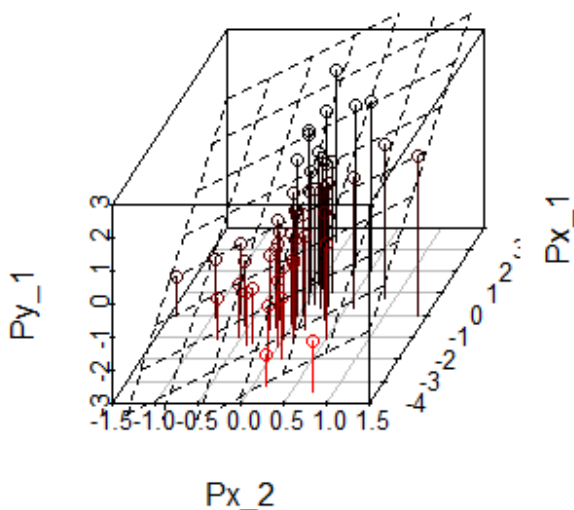
```
1 d1_1_scr<-data.frame(Px_2=Px1[,2],
2                       Px_1=Px1[,1],
3                       Py_1=Py1[,1])
4 d1_2_scr<-data.frame(Px_2=Px1[,2],
5                       Px_1=Px1[,1],
6                       Py_2=Py1[,2])
```



```

7
8 par(mfrow=c(1,2))
9
10 g3d1_1_scr<-scatterplot3d(d1_1_scr,
11     type="h",
12     angle=60,
13     highlight.3d = T)
14 fit1_1_scr<-lm(Py_1~Px_2+Px_1-1,d1_1_scr)
15 g3d1_1_scr$plane3d(fit1_1_scr)
16
17
18 g3d1_2_scr<-scatterplot3d(d1_2_scr,
19     type="h",
20     angle=60,
21     highlight.3d = T)
22 fit1_2_scr<-lm(Py_2~Px_2+Px_1-1,d1_2_scr)
23 g3d1_2_scr$plane3d(fit1_2_scr)

```



しかし、PLSはこうした主成分分析のスコアどうしのOLS回帰ではない！

PLS回帰

PLSの考え方

Y のスコアと X のスコアの内積 $(Xu^*)'Yv^*$ が最大化するように各スコアを決める。

Y のスコアと X のスコアの空間で主成分分析しているのと同じ考え方。

しかし、このスコアは X と Y をそれぞれ別々に主成分分析した場合のスコアとは違う。

$(Xu^*)'Yv^*$ が最大化するように u^* と v^* を決めた場合のスコア。

ただし、 $u'u = 1, v'v = 1$ に制約しておく。

式で表すと以下ということになる.

$$\max_{u,v} (Xu)'Yv$$

$$\text{s.t. } u'u = 1, v'v = 1$$

なので, ラグランジュ乗数法で

$$\max_{u,v} L = (Xu)'Yv - \lambda(u'u - 1) - \xi(v'v - 1)$$

最大化の一階の条件は, u で偏微分して...

$$X'Yv = 2\lambda u$$

v で偏微分して...

$$(X'Y)'u = 2\xi v$$

$$\lambda = \xi$$

であるなら, おお! これは特異値分解の必要条件そのものではないか!!!

...ということで, u, v は $X'Y$ の特異値分解から得られる?

$X'Y$ は...

```
1 xy1<-t(x1)%*%y1
2 round(xy1,4)
```

```
1 ##      [,1]  [,2]
2 ## [1,] 45.1963 30.8037
3 ## [2,] 48.9877  1.0123
```

$X'Y$ の特異値分解をすると...

```
1 svd_xy1<-svd(xy1)
2 print(svd_xy1)
```

```
1 ## $d
2 ## [1] 70.43287 20.77510
3 ##
4 ## $u
5 ##      [,1]  [,2]
6 ## [1,] -0.7518087 -0.6593812
7 ## [2,] -0.6593812  0.7518087
8 ##
9 ## $v
10 ##      [,1]  [,2]
11 ## [1,] -0.9410458  0.3382792
12 ## [2,] -0.3382792 -0.9410458
```

u^* と v^* が得られる...はず

```

1 u1<-svd_xy1$u
2 v1<-svd_xy1$v
3 d1<-svd_xy1$d

```

(符号が逆方向で気持ち悪いので、逆しておく)

```

1 u1<-u1
2 v1<-v1

```

ちなみに d は Xu と Yv の内積

```

1 t(x1*u1)*y1*v1>%
2 round(4)

```

```

1 ##      [,1]  [,2]
2 ## [1,] 70.4329  0.0000
3 ## [2,]  0.0000 20.7751

```

($n - 1 = 50$ で割ると共分散) .

$X'Yv = 2\lambda u$ の 2λ の値は？

d (つまり, Xu と Yv の内積) に一致する.

$X'Yv$ が . . .

```

1 t(x1)*y1*v1>%
2 round(4)

```

```

1 ##      [,1]  [,2]
2 ## [1,] 52.9520 13.6987
3 ## [2,] 46.4421 -15.6189

```

$2\lambda u$ が . . .

```

1 u1*diag(d1)

```

```

1 ##      [,1]  [,2]
2 ## [1,] 52.95205 13.69871
3 ## [2,] 46.44211 -15.61890

```

ぴったり.

$(X'Y)'u$ が . . .

```

1 t(t(x1)*y1)*u1

```

```

1 ##      [,1]  [,2]
2 ## [1,] 66.28055 -7.027784
3 ## [2,] 23.82597 19.550319

```

2ξが...

```
1 v1%*%diag(d1)
```

```
1 ##      [,1]      [,2]
2 ## [1,] 66.28055 -7.027784
3 ## [2,] 23.82597 19.550319
```

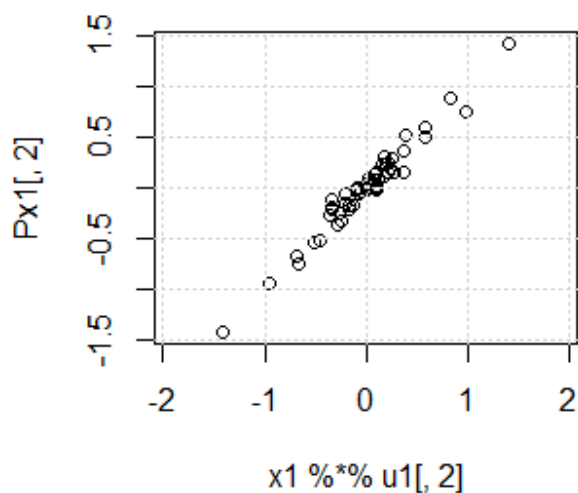
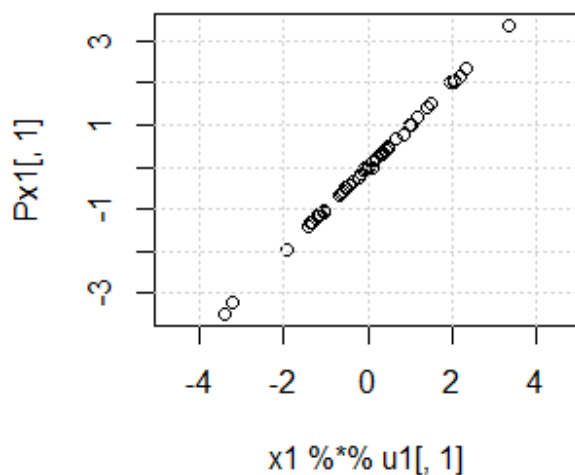
ぴったり。つまり、

$2\lambda = 2\xi = d$ ということでもある。

PCRとPLS

ノイズ（かく乱項）を入れていないので、 X と Y を、それぞれ主成分分析したスコアと、PLSの考え方で得られたスコアは一致するのか？

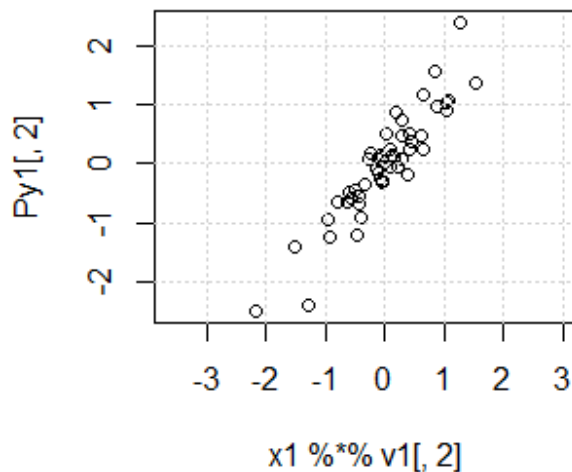
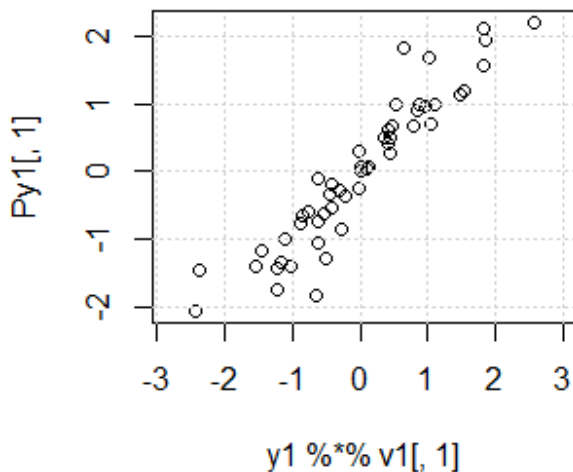
```
1 par(mfrow=c(1,2))
2 plot(Px1[,1]~x1%*%u1[,1],asp=1)
3 grid()
4 plot(Px1[,2]~x1%*%u1[,2],asp=1)
5 grid()
```



おおよそは一致しているが、**ぴったりとは一致していない**。

Y ではさらにぶれる。

```
1 par(mfrow=c(1,2))
2 plot(Py1[,1]~y1%*%v1[,1],asp=1)
3 grid()
4 plot(Py1[,2]~y1%*%v1[,2],asp=1)
5 grid()
```



PLSの考え方で最大化された $(Xu^*)'Yv^*$ が・・・

```
1 xyPx1<-x1%*%u1
2 xyPy1<-y1%*%v1
3 t(xyPx1)%*%xyPy1%>%
4 round(4)
```

```
1 ##      [,1]  [,2]
2 ## [1,] 70.4329  0.0000
3 ## [2,]  0.0000 20.7751
```

主成分分析で、 X と Y で個別に抽出した主成分の内積は・・・

```
1 t(Px1)%*%Py1%>%
2 round(4)
```

```
1 ##      Py1_1  Py1_2
2 ## PC1      63  31.1840
3 ## PC2      13 -16.7914
```

PLSの内積よりも小さい！

主成分分析では、 X と Y それぞれに、最も分散の大きな主成分を抽出しているはずだが、 X と Y との関係についての情報を含んでいないので、全体としては（PLSの意味での）最適ではないということ。

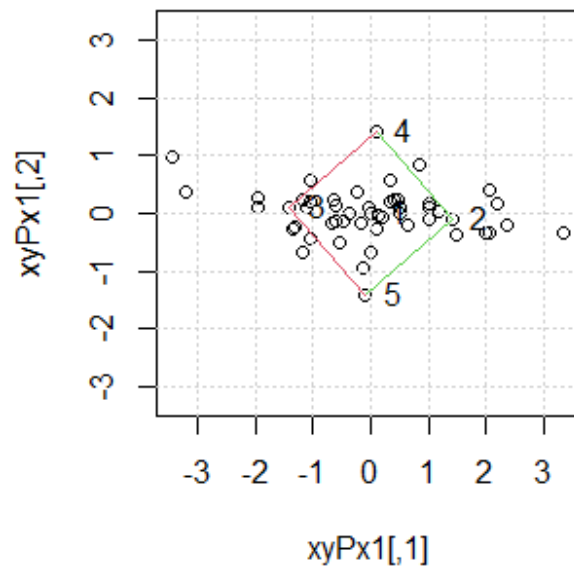
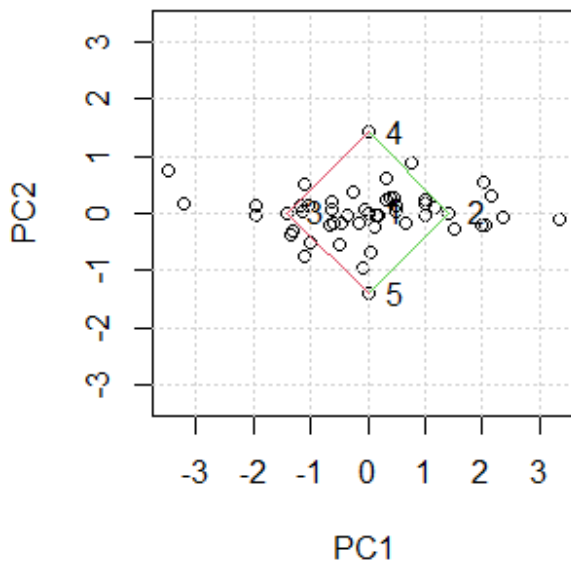
次のプロットは、左がPCRのスコア。右がPLSのスコア。

X のスコアはあまり変わらないが・・・

```

1 par(mfrow=c(1,2))
2 plot(Px1,asp=1)
3 grid()
4 text(x=Px1[1:5,1],y=Px1[1:5,2],1:5,pos=4)
5 f_grid1(Px1)
6
7 plot(xyPx1,asp=1)
8 grid()
9 text(x=xyPx1[1:5,1],y=xyPx1[1:5,2],1:5,pos=4)
10 f_grid1(xyPx1)

```

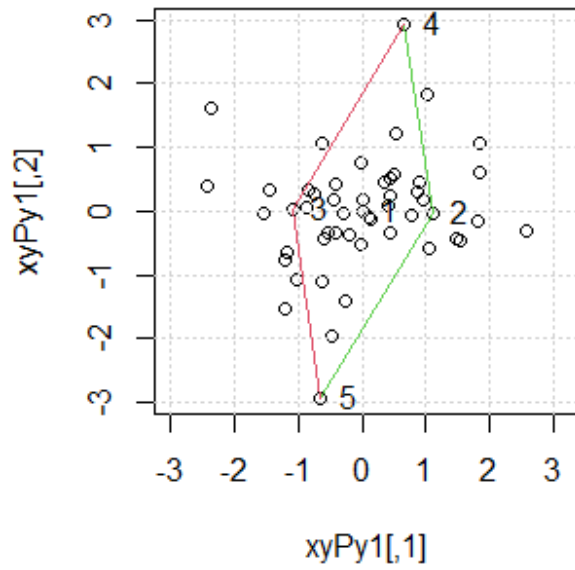
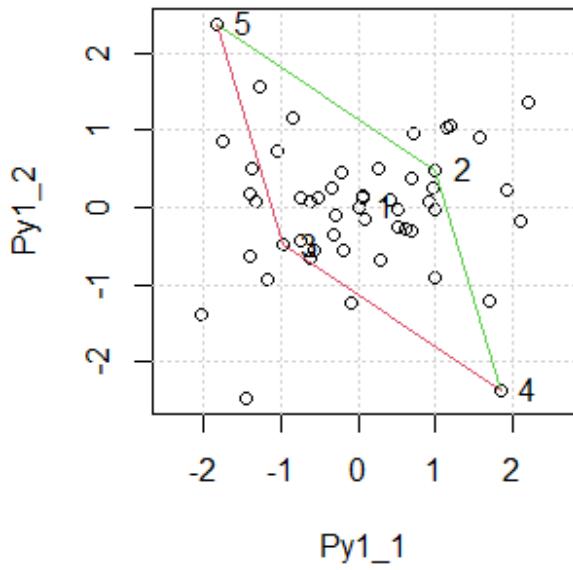


Yのスコアは、よりXに直交する形でとられているように見える。

```

1 par(mfrow=c(1,2))
2 plot(Py1,asp=1)
3 grid()
4 text(x=Py1[1:5,1],y=Py1[1:5,2],1:5,pos=4)
5 f_grid1(Py1)
6
7 plot(xyPy1,asp=1)
8 grid()
9 text(x=xyPy1[1:5,1],y=xyPy1[1:5,2],1:5,pos=4)
10 f_grid1(xyPy1)

```



$X'Y$ のスコアは直交していない

```
1 var(xyPx1)%>%
2 round(4)
```

```
1 ##      [,1]  [,2]
2 ## [1,]  1.7932 -0.1043
3 ## [2,] -0.1043  0.2068
```

```
1 var(xyPy1)%>%
2 round(4)
```

```
1 ##      [,1]  [,2]
2 ## [1,]  1.1401  0.1696
3 ## [2,]  0.1696  0.8599
```

スコアと X, Y との関係も、45度回転ではない。

```
1 xya1_x<-lm(xyPx1~x1-1)%>%
2 coef()
3 xya1_x
```

```
1 ##      [,1]  [,2]
2 ## x11  0.7518087  0.6593812
3 ## x12  0.6593812 -0.7518087
```

```
1 xya1_y<-lm(xyPy1~y1-1)%>%
2 coef()
3 xya1_y
```

```

1 ##          [,1]      [,2]
2 ## y11 0.9410458 -0.3382792
3 ## y12 0.3382792  0.9410458

```

スコアどうしの関係も定義した B とは異なる

```

1 xyb1<-lm(xyPy1~xyPx1-1)%>%
2   coef()
3 xyb1

```

```

1 ##          [,1]      [,2]
2 ## xyPx11 0.8093295 0.1204335
3 ## xyPx12 0.4083002 2.0696220

```

```

1 b0

```

```

1 ##          [,1]      [,2]
2 ## [1,]  0.7  0.346489
3 ## [2,]  1.3 -1.679139

```

X と Y の主成分どうしの関係はこちら側で定義して作成したのに・・・

$X \rightarrow X$ のスコア, $Y \rightarrow Y$ のスコア, $X \rightarrow Y$ のスコア

もOLSで評価すると間違いないのに, PLSで一緒に分析すると, (PLSの意味で) もっとよい解が見つかった!・・・ということ.

ただし, PLSの係数を使って Y に直接 X を回帰した X の係数にあたる ABA^{-1} を計算すると,・・・

```

1 xya1_x%*%xyb1%*%solve(xya1_y)

```

```

1 ##          y11      y12
2 ## x11 0.3336750  1.666325
3 ## x12 0.7128141 -1.312814

```

(途中が違ってても) PCRの場合と (OLSとも) 一致する.

```

1 lm(y1~x1-1)%>%
2   coef()

```

```

1 ##          [,1]      [,2]
2 ## x11 0.3336750  1.666325
3 ## x12 0.7128141 -1.312814

```

(ノイズが一切ないので当然だが)

Rの `pls` パッケージを試してみる.

データ

データは、data.frameにするが、 X と Y をそれぞれmatrixとして、2変数data.frameにするという荒業を使う。

matrixの `x1` と `y1` をそのままdata.frameに突っ込むには、`I()` 関数（そのまんまにしてね！という関数）を使う。

それぞれ `x` と `y` と名前を付けよう。

```
1 d01<-data.frame(x=I(x1),y=I(y1))
```

データ構造を見てみると・・・

```
1 str(d01)
```

```
1 ## 'data.frame': 51 obs. of 2 variables:
2 ## $ x: 'AsIs' num [1:51, 1:2] 0 1 -1 1 -1 ...
3 ## $ y: 'AsIs' num [1:51, 1:2] -1.06e-18 1.05 -1.05 -3.79e-01 3.79e-01 ...
```

51行の2変数で、`x` と `y` があり、それぞれ、`[1:51,1:2]` の変数であると説明にある。

pls回帰

回帰には、`pls()` 関数を使う。引数は、`formula=` の後に・・・

`ncomp=` は抽出する主成分の数（この場合、2変数しかないなので、それ以上は無理）

`data=` はさっきつくったdata.frame

`validation=` は、クロスバリデーション（CV）の方法。“LOO”がLeave-one-out法。この場合、ややこしいので `=“none”`（CV）にした。

```
1 fit01<-pls(y~x,
2           ncomp=2,
3           data=d01,
4           validation="none")
```

要約を表示しても、主成分の寄与率しか表示してくれない。

```
1 summary(fit01)
```

```
1 ## Data: X dimension: 51 2
2 ## Y dimension: 51 2
3 ## Fit method: kernelpls
4 ## Number of components considered: 2
5 ## TRAINING: % variance explained
6 ## 1 comps 2 comps
7 ## X 89.96 100
8 ## Y1 98.00 100
9 ## Y2 12.66 100
```

PLS回帰を使う人が、回帰係数よりも、たくさんの説明変数の中で**どの変数との関係が強い**かに関心がある場合に使われているということか？

適合度とモデル選択

適合度は、Root Mean Squared Error of Prediction (RMSEP) で表す。

なんか長いが、要は、実測値 y_i と予測値 \hat{y}_i との差の2乗平均の平方根をとったもの。当然小さい方がよいですね。

$$RMSEP = \sqrt{\frac{1}{n} \sum_i (y_i - \hat{y})^2}$$

`RMSEP` 関数が用意されている。

```
1 RMSEP(fit01)
```

```
1 ##
2 ## Response: Y1
3 ## (Intercept)      1 comps      2 comps
4 ##  9.901e-01    1.401e-01    1.871e-16
5 ##
6 ## Response: Y2
7 ## (Intercept)      1 comps      2 comps
8 ##  9.901e-01    9.253e-01    8.450e-16
```

`(Intercept)` は、何も説明変数を使わなかった Y のRMSEP.ほぼ標準偏差ということになるが、 $n - 1$ ではなく、 n で割られるので、1となっていない。

```
1 sd(y1[,1])*sqrt(50/51)
```

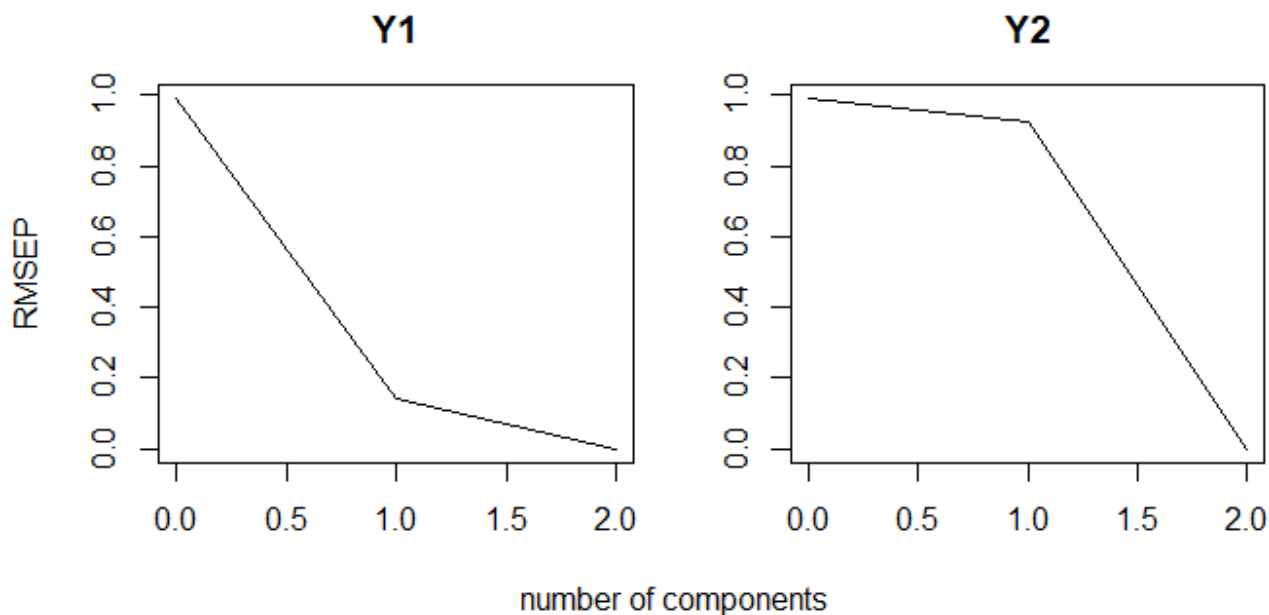
```
1 ## [1] 0.9901475
```

`1 comps` は、第1主成分だけを使った場合のRMSEP.

`2 comps` は、第2主成分まで使った場合のRMSEP. 2変数しかないので、当然ここで0となる。

プロットもできる。

```
1 plot(RMSEP(fit01))
```



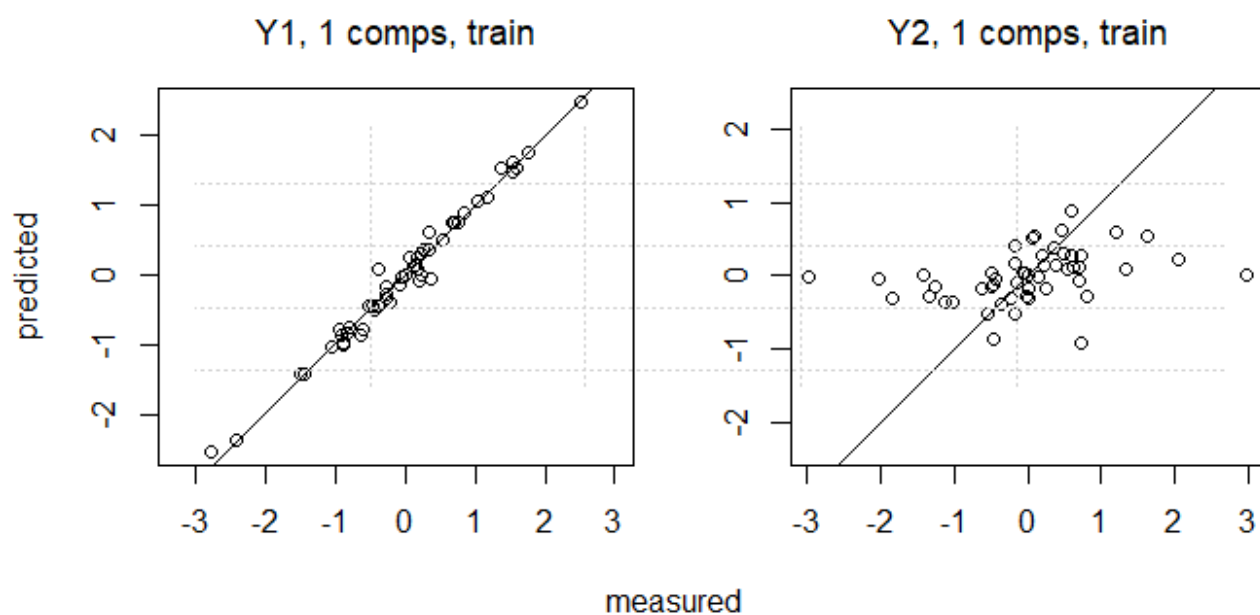
フィットの様子を確認することもできる。

横軸は抽出した主成分の数だが、1.2とか、余計な目盛がついているので、それは無視。主成分名を数字に置き換えてあるので、主成分の数が少ないところなるみたいだ。

そのまま、`plot()` 関数を使うと、Yの実際値と予測値との相関を出力する。

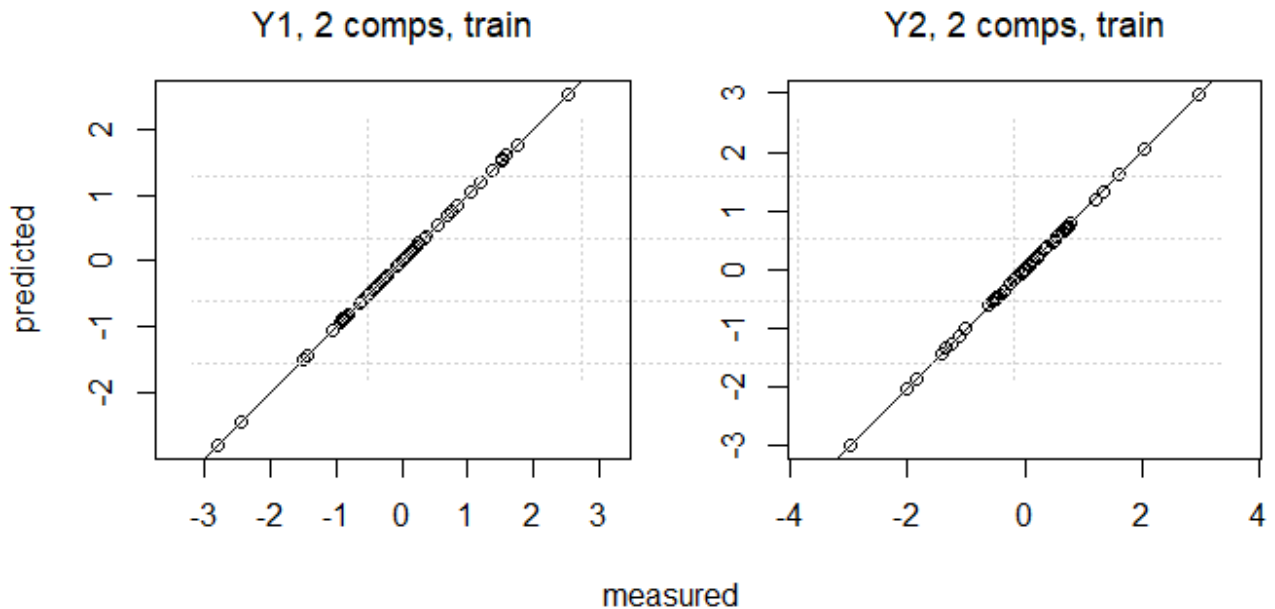
第1主成分だけで予測した場合

```
1 plot(fit01,ncomp=1,line=TRUE,asp=1)
2 grid()
```



第2主成分も使って予測した場合

```
1 plot(fit01,ncomp=2,line=TRUE,asp=1)
2 grid()
```



ちなみに、予測値は `predict` 関数で出力できる。

```
1 predict(fit01,comps=1)%>%
2   round(.,4)%>%
3   head()
```

```
1 ##      Y1      Y2
2 ## 1  0.0000  0.0000
3 ## 2  1.0432  0.3750
4 ## 3 -1.0432 -0.3750
5 ## 4  0.0683  0.0246
6 ## 5 -0.0683 -0.0246
7 ## 6 -0.7659 -0.2753
```

主成分のスコア

X のスコアは `scores()` 関数で

```
1 Scrx1<-scores(fit01)
2 Scrx1%>%
3   round(4)%>%
4   head()
```

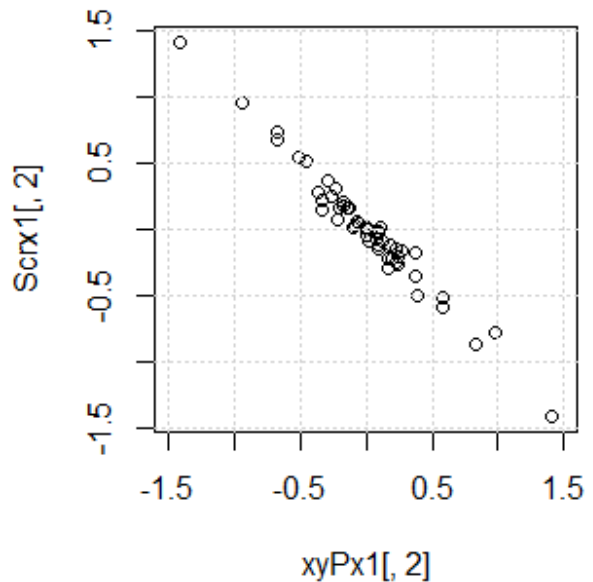
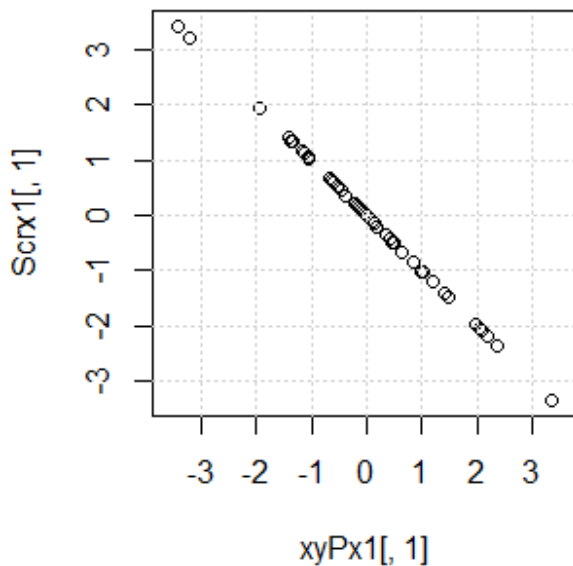
```
1 ##      Comp 1  Comp 2
2 ## 1  0.0000  0.0000
3 ## 2 -1.4112  0.0103
4 ## 3  1.4112 -0.0103
5 ## 4 -0.0924 -1.4166
6 ## 5  0.0924  1.4166
7 ## 6  1.0361 -0.1252
```

X のスコアは Xu^* から出した `xyScrx1` に（符号を除いて）一致するはずだが、

```

1 par(mfrow=c(1,2))
2 plot(Scrx1[,1]~xyPx1[,1],asp=1)
3 grid()
4 plot(Scrx1[,2]~xyPx1[,2],asp=1)
5 grid()

```



第2主成分で多少の違いが・・・

Yのスコアは `Yscores()` 関数で

```

1 Scry1<-Yscores(fit01)
2 Scry1%>%
3   round(4)%>%
4   head()

```

```

1 ##   Comp 1  Comp 2
2 ## 1  0.0000  0.0000
3 ## 2 -1.4058  0.0103
4 ## 3  1.4058 -0.0103
5 ## 4 -0.8287 -1.4166
6 ## 5  0.8287  1.4166
7 ## 6  0.9710 -0.1252

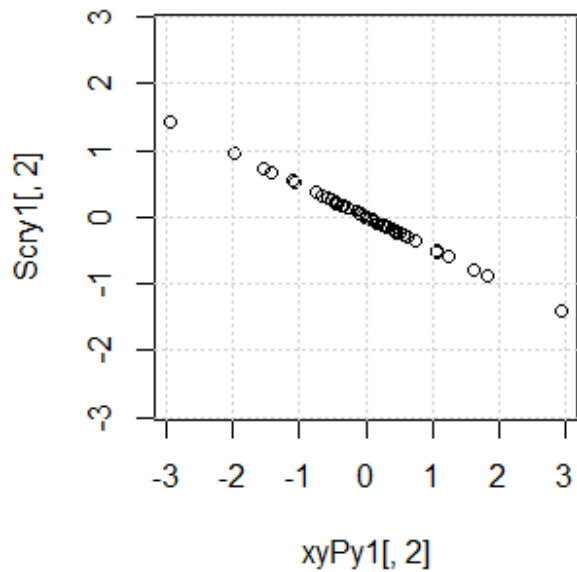
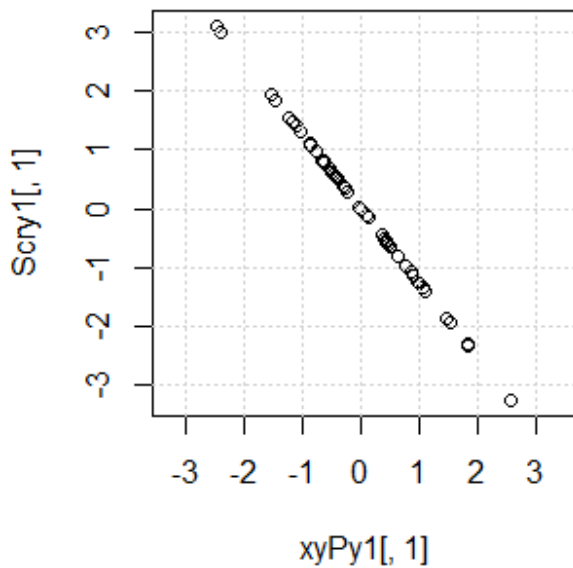
```

Yについては、比例はするが、スケールが全然ちがう。なぜ???

```

1 par(mfrow=c(1,2))
2 plot(Scry1[,1]~xyPy1[,1],asp=1)
3 grid()
4 plot(Scry1[,2]~xyPy1[,2],asp=1)
5 grid()

```



`plsrf()` がYのスコアのばらつきをXのスコアのばらつきに対応させている？

```
1 lm(Yscores(fit01)~scores(fit01)-1)%>%
2   coef()%>%
3   round(4)
```

```
1 ##           Comp 1 Comp 2
2 ## scores(fit01)Comp 1 1.0000    0
3 ## scores(fit01)Comp 2 0.5198    1
```

Yのスコアの分散は

```
1 var(Scry1)%>%
2   round(4)
```

```
1 ##           Comp 1 Comp 2
2 ## Comp 1 1.8474 0.1043
3 ## Comp 2 0.1043 0.2008
```

X'Yから求めたYのスコアの分散は

```
1 var(xyPy1)%>%
2   round(4)
```

```
1 ##           [,1] [,2]
2 ## [1,] 1.1401 0.1696
3 ## [2,] 0.1696 0.8599
```

`plsrf` で求めたYの第1主成分のスコアは、約1.27倍されている。

```
1 sqrt(1.8474/1.1401)
```

```
1 ## [1] 1.272943
```

第2主成分のスコアは約0.48倍されている

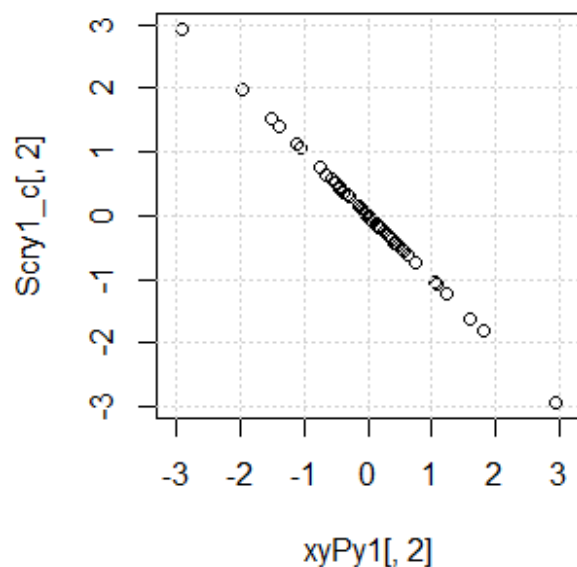
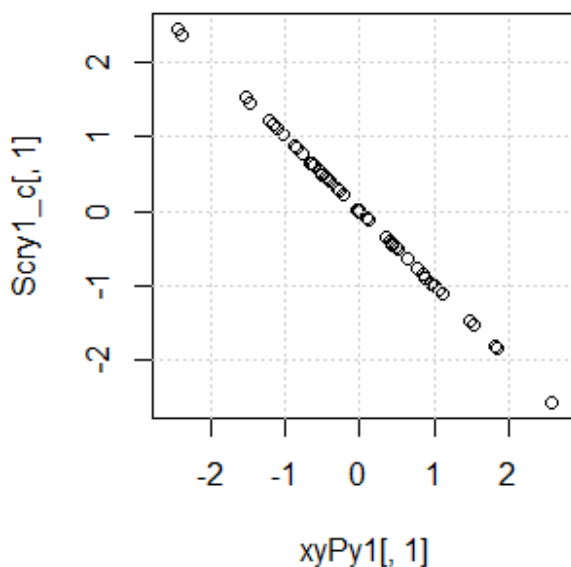
```
1 sqrt(0.2008/0.8599)
```

```
1 ## [1] 0.4832344
```

`plsr` の2つのスコアを、それぞれ割ってやると、 X 、 Y の最大化から求めた Y のスコアとほぼ一致する。

```
1 Scry1_c<-sweep(Scry1,2,c(sqrt(1.8474/1.1401),sqrt(0.2008/0.8599)),FUN="/")
```

```
1 par(mfrow=c(1,2))
2 plot(Scry1_c[,1]~xyPy1[,1],asp=1)
3 grid()
4 plot(Scry1_c[,2]~xyPy1[,2],asp=1)
5 grid()
```

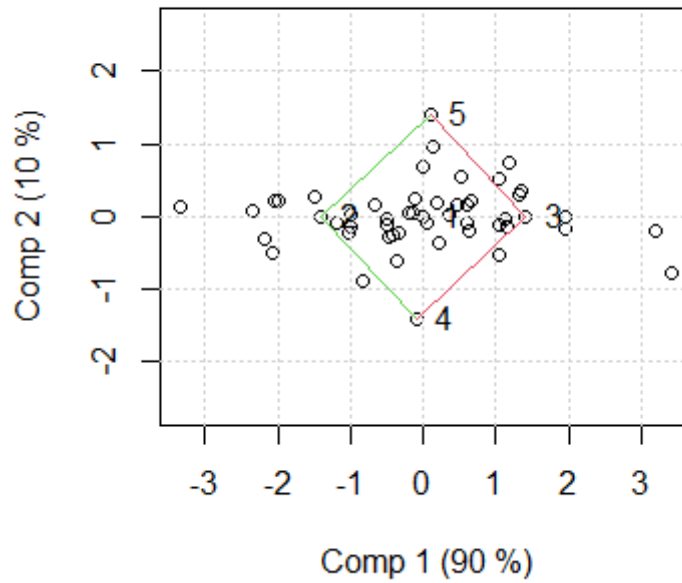


なんでそんなことするかは謎？

Xのスコアのプロット

`plot` 関数に引数 `plottype="scores"` で、`scoreplot()` 関数でもOK

```
1 plot(fit01,plottype="scores",comps=1:2,asp=1)
2 grid()
3 text(x=Scrx1[1:5,1],y=Scrx1[1:5,2],1:5,pos=4)
4 f_grid1(Scrx1)
```



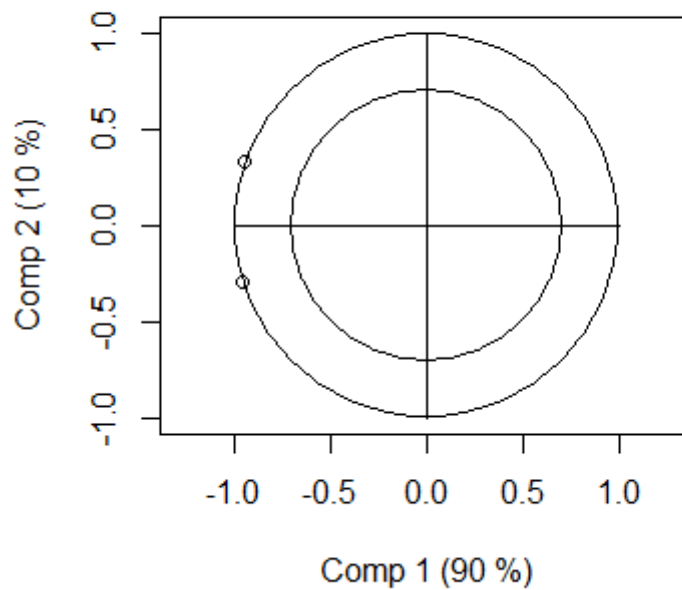
多変数の場合は、`pairs()` の図が出力されるようだ・・・

Yのスコアは関数が用意されていないが、この場合Xのスコアと全く同じなので、あんまり意味ない。

スコアとXとの相関のプロット

`plot()` 関数の `plottype="correlation"` 引数でも同じ。

```
1 | corplot(fit01)
```




```
1 cor(x1,Scrx1)%>%
2 round(4)
```

```
1 ##      Comp 1  Comp 2
2 ## [1,] -0.9554 -0.2954
3 ## [2,] -0.9416  0.3369
```

・・・というのをプロットしているようだ。左下が x_1 で、左上が x_2 。

x_1 も x_2 も、第1主成分の方に相関が高いようだ・・・ということはわかる。

PLSのOLSやPCRとの違い

以上は全くノイズの無い世界の話。ノイズがなければ、OLSで X を直接 Y に回帰させても、あるいは、主成分回帰(PCR)で、いったん X の主成分を抽出して、主成分を Y に回帰させても、それはどう解釈するかが違うだけで、情報の損失はない。

しかし、 X は2変数よりもすごく多いと、そうはいかなくなる。

まず、OLSだと X の変数間でマルチコが生じる可能性が高まる。

だからといって、PCRを使うとすると、 X から主成分を抽出して、**主な**主成分を2個とか3個とか選んで使うことになる。

ところが、この主成分は、あくまでも X の中だけの情報で抽出されるし、どれがどれが**主な**主成分かも X の中だけで決まる。これには Y の情報は入らない。

しかし、PLSだと Y の情報も入れて主成分を抽出するので、PCRよりもうまく主成分を抽出できる・・・という話。

以下、ノイズを入れてやってみよう。

2×2変数にノイズあり

データ作成

攪乱項を加える

2変数×2変数というのは変えないで、上記で作成した `x1` や `y1` にちょっとだけかく乱項を入れてみる。

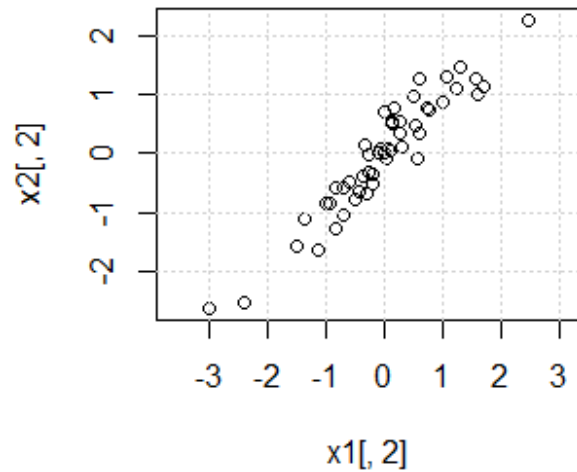
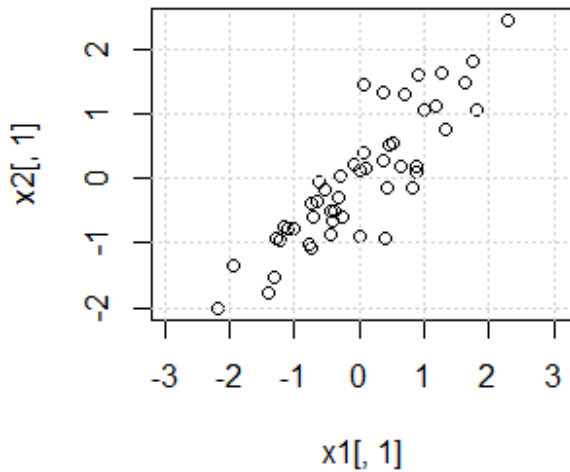
```
1 x2<-x1
2 set.seed(111)
3 x2[6:51,1]<-x1[6:51,1]+rnorm(46,mean=0,sd=0.5)
4 set.seed(124)
5 x2[6:51,2]<-x1[6:51,2]+rnorm(46,mean=0,sd=0.5)
6 x2<-apply(x2,2,scale)
```

X にどのくらいノイズが入ったかプロットしてみる

```

1 par(mfrow=c(1,2))
2 plot(x2[,1]~x1[,1],asp=1)
3 grid()
4 plot(x2[,2]~x1[,2],asp=1)
5 grid()

```



このぐらい・・・

Yにもノイズを加える

```

1 y2<-y1
2 set.seed(125)
3 y2[6:51,1]<-y2[6:51,1]+rnorm(46,0,0.5)
4 set.seed(126)
5 y2[6:51,2]<-y2[6:51,2]+rnorm(46,0,0.5)

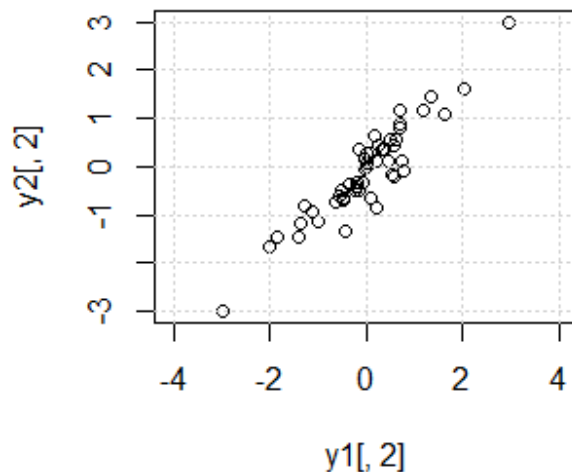
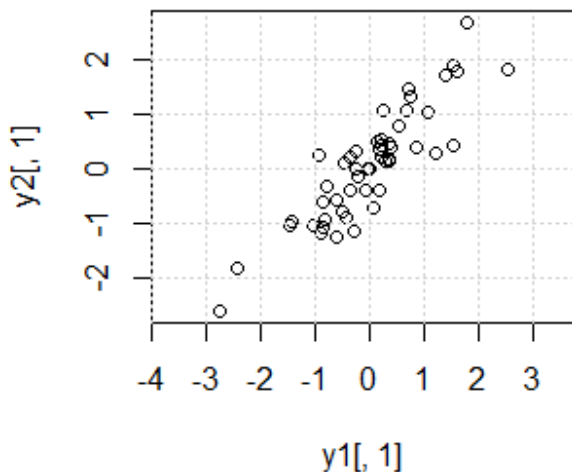
```

Yにどれだけノイズが入ったか見てみよう

```

1 par(mfrow=c(1,2))
2 plot(y2[,1]~y1[,1],asp=1)
3 grid()
4 plot(y2[,2]~y1[,2],asp=1)
5 grid()

```



このぐらい.

OLSやPCRとの比較

OLS回帰

作成した X と Y で OLS 回帰してみる.

```
1 lm(y2~x2-1)%>%
2   coef()
```

```
1 ##           [,1]      [,2]
2 ## x21  0.3545037  0.9186370
3 ## x22  0.5574486 -0.6727718
```

真の値はこれ・・・

```
1 a1%*%b0%*%solve(a1)
```

```
1 ##           [,1]      [,2]
2 ## [1,]  0.3336750  1.666325
3 ## [2,]  0.7128141 -1.312814
```

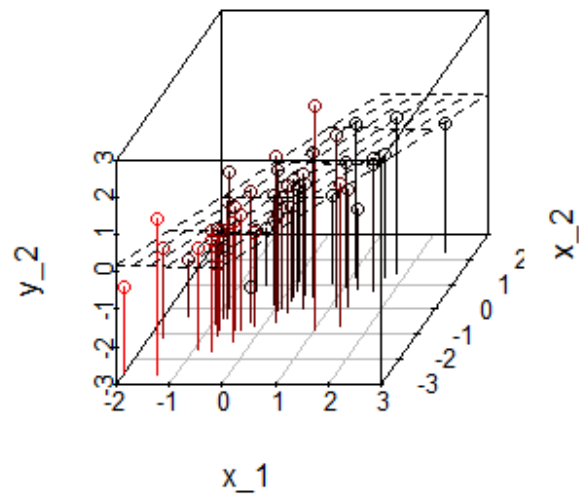
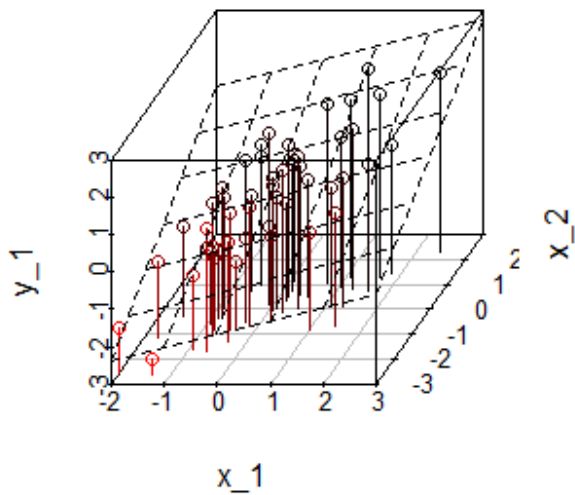
y_1 の x_1 の係数以外はあてはまりがいまいち。こういうのは、乱数を多数発生させて、推定値の分布を見て判断すべきで、1つの例でどうのこうの言えることではないのだが、 y_1 をプロットしてみると・・・

```
1 d2_1_OLS<-data.frame(x_1=x2[,1],
2                     x_2=x2[,2],
3                     y_1=y2[,1])
4 d2_2_OLS<-data.frame(x_1=x2[,1],
5                     x_2=x2[,2],
6                     y_2=y2[,2])
7 par(mfrow=c(1,2))
8 g3d1_1_OLS<-scatterplot3d(d2_1_OLS,
```

```

9         type="h",
10        angle=60,
11        highlight.3d = T)
12 fit2_1_OLS<-lm(y_1~x_1+x_2-1,d2_1_OLS)
13 g3d1_1_OLS$plane3d(fit2_1_OLS)
14
15 g3d1_2_OLS<-scatterplot3d(d2_2_OLS,
16        type="h",
17        angle=60,
18        highlight.3d = T)
19 fit2_2_OLS<-lm(y_2~x_1+x_2-1,d2_2_OLS)
20 g3d1_2_OLS$plane3d(fit2_2_OLS)

```



Xの相関に対して回帰平面が不安定そうなのは確認できる.

係数のt値では, そんな不安定であるようには見えないのが問題.

```

1 lm(y2~x2-1)%>%
2   summary()%>%
3   coefficients()

```

```

1 ## Response Y1 :
2 ##      Estimate Std. Error t value    Pr(>|t|)
3 ## x21  0.3545037  0.1118888  3.168358 2.640143e-03
4 ## x22  0.5574486  0.1118888  4.982167 8.223860e-06
5 ##
6 ## Response Y2 :
7 ##      Estimate Std. Error t value    Pr(>|t|)
8 ## x21  0.9186370  0.1255169  7.318831 2.130580e-09
9 ## x22 -0.6727718  0.1255169 -5.360010 2.226719e-06

```

PCR 回帰

X を主成分分析

```
1 prc_x2<-prcomp(x2)
2 prc_x2
```

```
1 ## Standard deviations (1, ..., p=2):
2 ## [1] 1.2830922 0.5947054
3 ##
4 ## Rotation (n x k) = (2 x 2):
5 ##           PC1      PC2
6 ## [1,] -0.7071068  0.7071068
7 ## [2,] -0.7071068 -0.7071068
```

主成分分析のスコア（方向を逆に）

```
1 Px2<-prc_x2$x
2 Px2[,1]<-prc_x2$x[,1]
```

Y を X のスコアで回帰

```
1 lm(y2~Px2-1)%>%
2   coef()%>%
3   round(4)
```

```
1 ##           [,1]  [,2]
2 ## Px2PC1  0.6448 0.1739
3 ## Px2PC2 -0.1435 1.1253
```

真の値は・・・

```
1 b0%*%solve(a1)
```

```
1 ##           [,1]  [,2]
2 ## [1,]  0.7399795 0.249970
3 ## [2,] -0.2680918 2.106569
```

これも1つの乱数だけで云々できないが、 y_2 に対する係数がだいぶ外れているようだ。

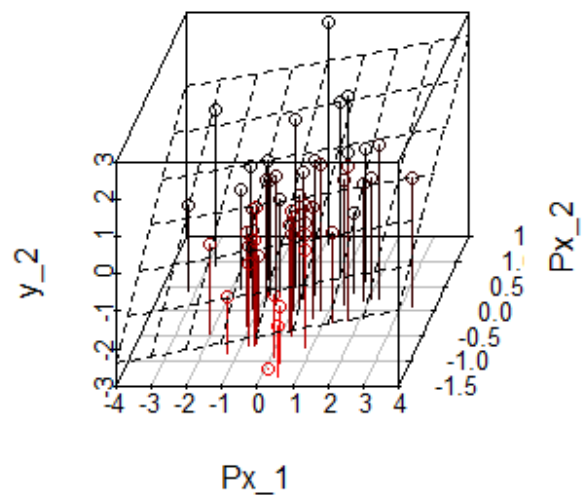
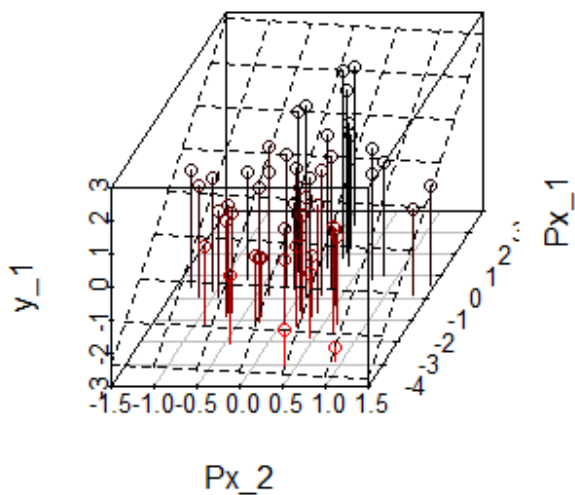
プロットしてみると・・・

```
1 d2_1_PCR<-data.frame(Px_2=Px2[,2],
2                       Px_1=Px2[,1],
3                       y_1=y2[,1])
4 d2_2_PCR<-data.frame(Px_1=Px2[,1],
5                       Px_2=Px2[,2],
6                       y_2=y2[,2])
7 par(mfrow=c(1,2))
8
9 g3d1_1_PCR<-scatterplot3d(d2_1_PCR,
10                          type="h",
11                          angle=60,
```

```

12 highlight.3d = T)
13 fit2_1_PCR<-lm(y_1~Px_2+Px_1-1,d2_1_PCR)
14 g3d1_1_PCR$plane3d(fit2_1_PCR)
15
16 g3d1_2_PCR<-scatterplot3d(d2_2_PCR,
17 type="h",
18 angle=60,
19 highlight.3d = T)
20 fit2_2_PCR<-lm(y_2~Px_1+Px_2-1,d2_2_PCR)
21 g3d1_2_PCR$plane3d(fit2_2_PCR)

```



説明変数の相関がましになったことは確認できる。

```

1 lm(y1~Px2-1)%>%
2 summary()%>%
3 coefficients()

```

```

1 ## Response Y1 :
2 ##           Estimate Std. Error  t value    Pr(>|t|)
3 ## Px2PC1  0.7360053  0.03174716  23.183344  4.577494e-28
4 ## Px2PC2 -0.2756515  0.06849531  -4.024385  1.975324e-04
5 ##
6 ## Response Y2 :
7 ##           Estimate Std. Error  t value    Pr(>|t|)
8 ## Px2PC1  0.2334963  0.07711698  3.027820  3.921328e-03
9 ## Px2PC2  1.1032709  0.16638187  6.630956  2.475341e-08

```

PLS回帰

主成分の抽出

$X'Y$ の特異値分解

```

1 svd_xy2<-svd(t(x2)%*%y2)
2 svd_xy2

```

```

1  ## $d
2  ## [1] 55.05462 19.84587
3  ##
4  ## $u
5  ##           [,1]      [,2]
6  ## [1,] -0.7462628 -0.6656515
7  ## [2,] -0.6656515  0.7462628
8  ##
9  ## $v
10 ##           [,1]      [,2]
11 ## [1,] -0.9599648  0.2801207
12 ## [2,] -0.2801207 -0.9599648

```

```

1  d2<-svd_xy2$d
2  u2<-svd_xy2$u
3  v2<-svd_xy2$v

```

XとYの主成分

```

1  scr_x2<-x2%*%u2
2  scr_y2<-y2%*%v2

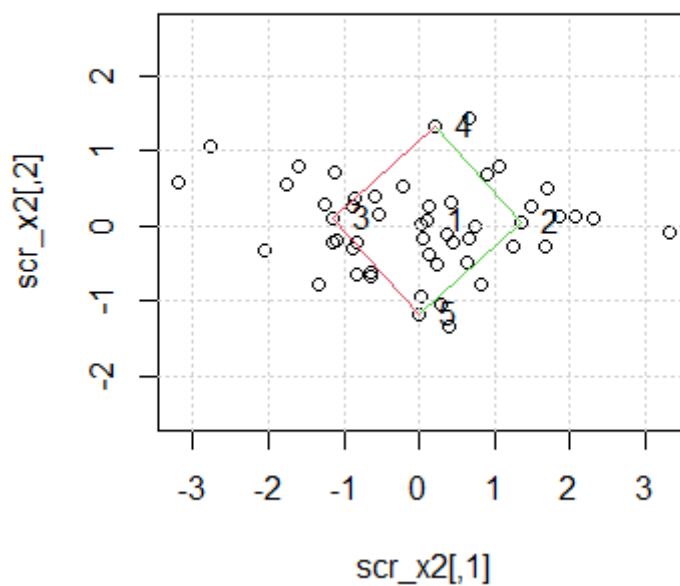
```

Xの主成分のスコア

```

1  plot(scr_x2,asp=1)
2  grid()
3  text(x=scr_x2[1:5,1],y=scr_x2[1:5,2],1:5,pos=4)
4  f_grid1(scr_x2)

```

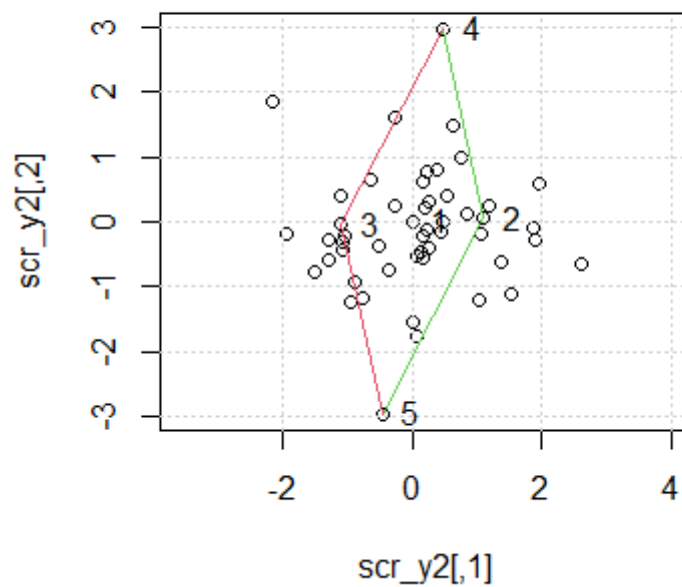


Yの主成分のスコア

```

1 plot(scr_y2,asp=1)
2 grid()
3 text(x=scr_y2[1:5,1],y=scr_y2[1:5,2],1:5,pos=4)
4 f_grid1(scr_y2)

```



ちなみに、 X と Y 、別々に主成分分析して主成分を抽出すると...

```

1 prc_x2<-prcomp(x2)
2 prc_y2<-prcomp(y2)

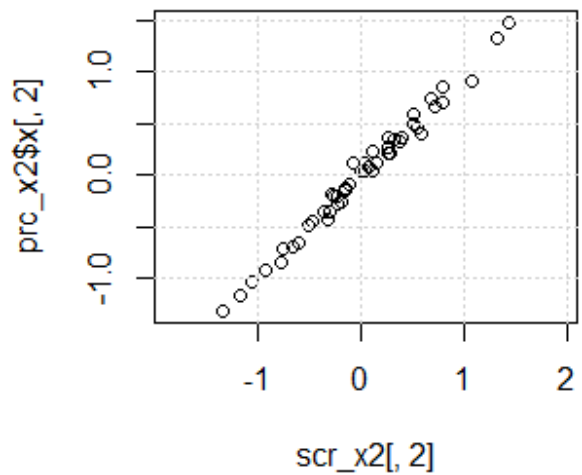
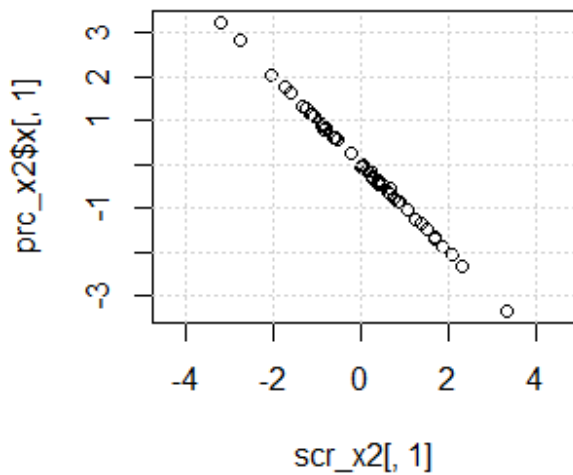
```

X ・ Y からの主成分と比較すると、

```

1 par(mfrow=c(1,2))
2 plot(prc_x2$x[,1]~scr_x2[,1],asp=1)
3 grid()
4 plot(prc_x2$x[,2]~scr_x2[,2],asp=1)
5 grid()

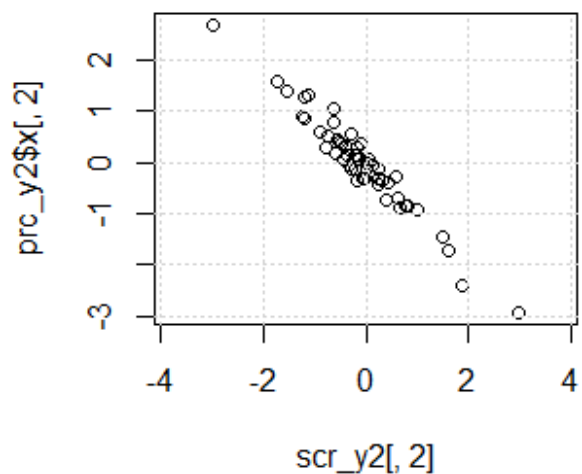
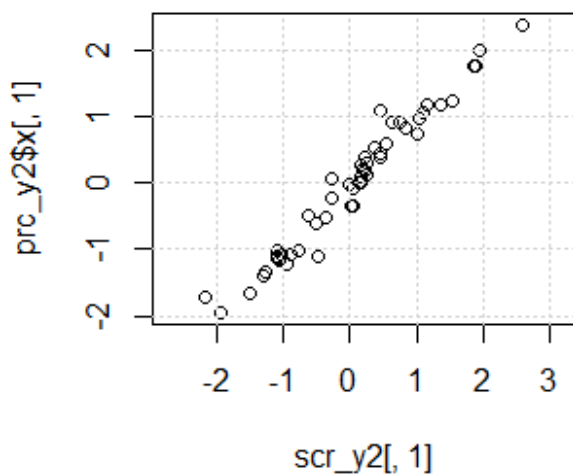
```

Xの主成分分析は若干ずれていることがわかる。

Yの主成分については、もっとずれてくる。

```
1 par(mfrow=c(1,2))
2 plot(prc_y2$x[,1]~scr_y2[,1],asp=1)
3 grid()
4 plot(prc_y2$x[,2]~scr_y2[,2],asp=1)
5 grid()
```



スコアどうして回帰する

```
1 lm(scr_y2~scr_x2-1)%>%
2 coef()
```

```
1 ##           [,1]      [,2]
2 ## scr_x21 0.6767605 0.05014599
3 ## scr_x22 0.1391105 1.11940449
```

真値は

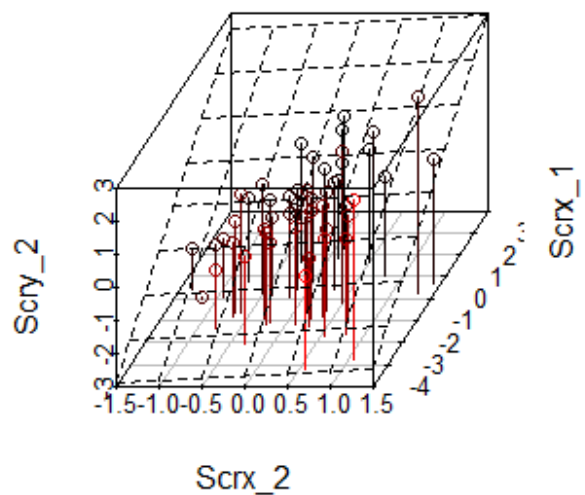
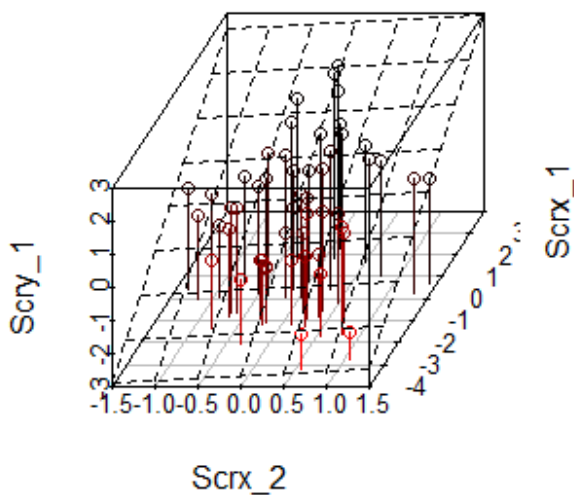
```
1 xyb1
```

```
1 ##           [,1]      [,2]
2 ## xyPx11 0.8093295 0.1204335
3 ## xyPx12 0.4083002 2.0696220
```

第2主成分が今一つだが、第1主成分は安定している。

プロットしてみると・・・

```
1 d2_1_PLS<-data.frame(Scrx_2=scr_x2[,2],
2                       Scrx_1=scr_x2[,1],
3                       Scry_1=scr_y2[,1])
4 d2_2_PLS<-data.frame(Scrx_2=scr_x2[,2],
5                       Scrx_1=scr_x2[,1],
6                       Scry_2=scr_y2[,2])
7
8 par(mfrow=c(1,2))
9
10 g3d2_1_PLS<-scatterplot3d(d2_1_PLS,
11                           type="h",
12                           angle=60,
13                           highlight.3d = T)
14 fit2_1_PLS<-lm(Scry_1~Scrx_2+Scrx_1-1,d2_1_PLS)
15 g3d2_1_PLS$plane3d(fit2_1_PLS)
16
17 g3d2_2_PLS<-scatterplot3d(d2_2_PLS,
18                           type="h",
19                           angle=60,
20                           highlight.3d = T)
21 fit2_2_PLS<-lm(Scry_2~Scrx_2+Scrx_1-1,d2_2_PLS)
22 g3d2_2_PLS$plane3d(fit2_1_PLS)
```



被説明変数がある程度ばらついており、 Y の傾向もPCRより確認できやすい気がする？

```
1 lm(scr_y2~scr_x2-1)%>%
2   summary()%>%
3   coefficients()
```

```
1 ## Response Y1 :
2 ##           Estimate Std. Error  t value    Pr(>|t|)
3 ## scr_x21  0.6767605  0.06498591  10.4139573 5.147398e-14
4 ## scr_x22  0.1391105  0.13920572   0.9993159 3.225513e-01
5 ##
6 ## Response Y2 :
7 ##           Estimate Std. Error  t value    Pr(>|t|)
8 ## scr_x21  0.05014599  0.07678037   0.6531095 5.167386e-01
9 ## scr_x22  1.11940449  0.16447051   6.8061106 1.324866e-08
```

Y の第1主成分について、 X の第1主成分がだけが効いており、 Y の第2主成分について、 X の第2主成分だけが効いている。

Y の主成分が2つ X の主成分が2つの場合、内積を最大化しようとする、主成分1対1で対応するよな気がする（要確認）。

Rの `p1sr()` 関数で

データ作成

```
1 d02<-data.frame(x=I(x2),y=I(y2))
```

PLS回帰

```
1 fit02<-plsr(y~x,  
2           ncomp=2,  
3           data=d02,  
4           validation = "none")
```

要約

```
1 summary(fit02)
```

```
1 ## Data:   X dimension: 51 2  
2 ## Y dimension: 51 2  
3 ## Fit method: kernelpls  
4 ## Number of components considered: 2  
5 ## TRAINING: % variance explained  
6 ##      1 comps  2 comps  
7 ## X      82.271  100.00  
8 ## Y1     65.185   66.29  
9 ## Y2      6.181   53.09
```

適合度とモデル選択

$$RMSEP = \sqrt{\frac{1}{n} \sum_i (y_i - \hat{y})^2}$$

```
1 RMSEP(fit02)
```

```
1 ##  
2 ## Response: Y1  
3 ## (Intercept)      1 comps      2 comps  
4 ##      1.0116      0.5969      0.5874  
5 ##  
6 ## Response: Y2  
7 ## (Intercept)      1 comps      2 comps  
8 ##      0.9586      0.9285      0.6565
```

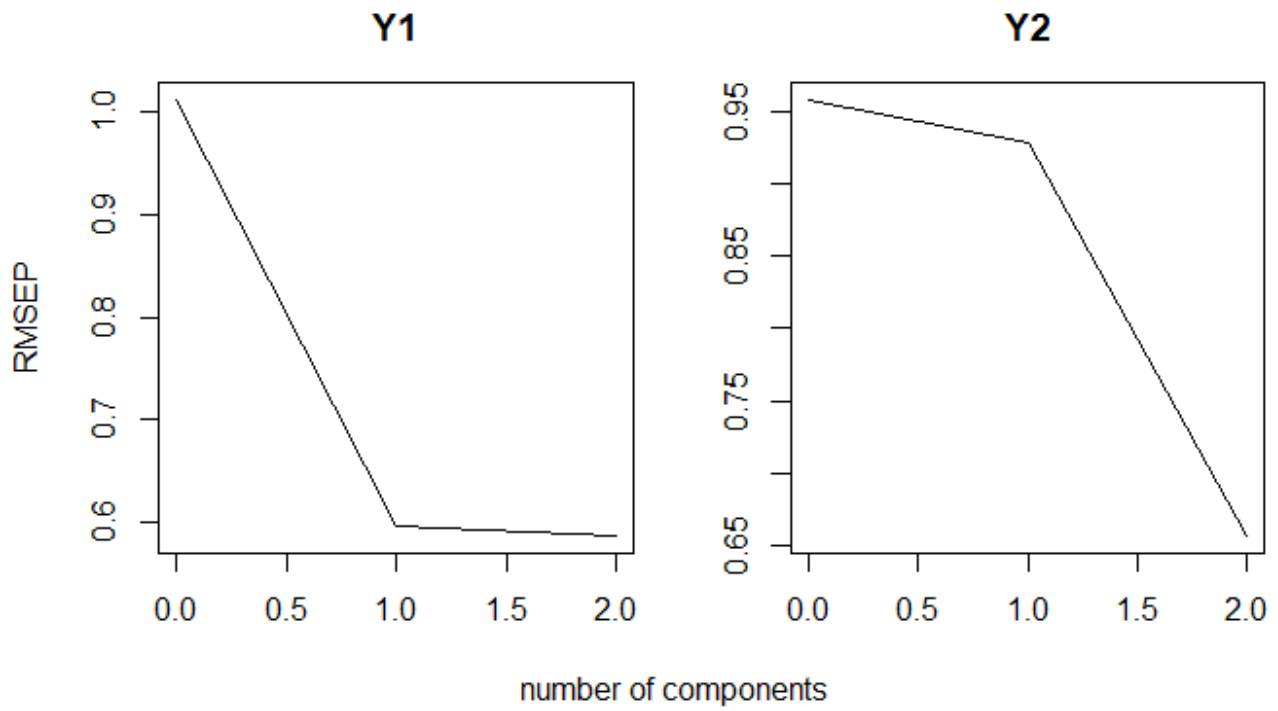
(Intercept) は、何もXを使わなかったときのRMSEP.

1 comps は、第1主成分だけを使った場合のRMSEP.

2 comps は、第2主成分まで使った場合のRMSEP. 2変数しかないが、ノイズがあるので0にはならない.

プロット

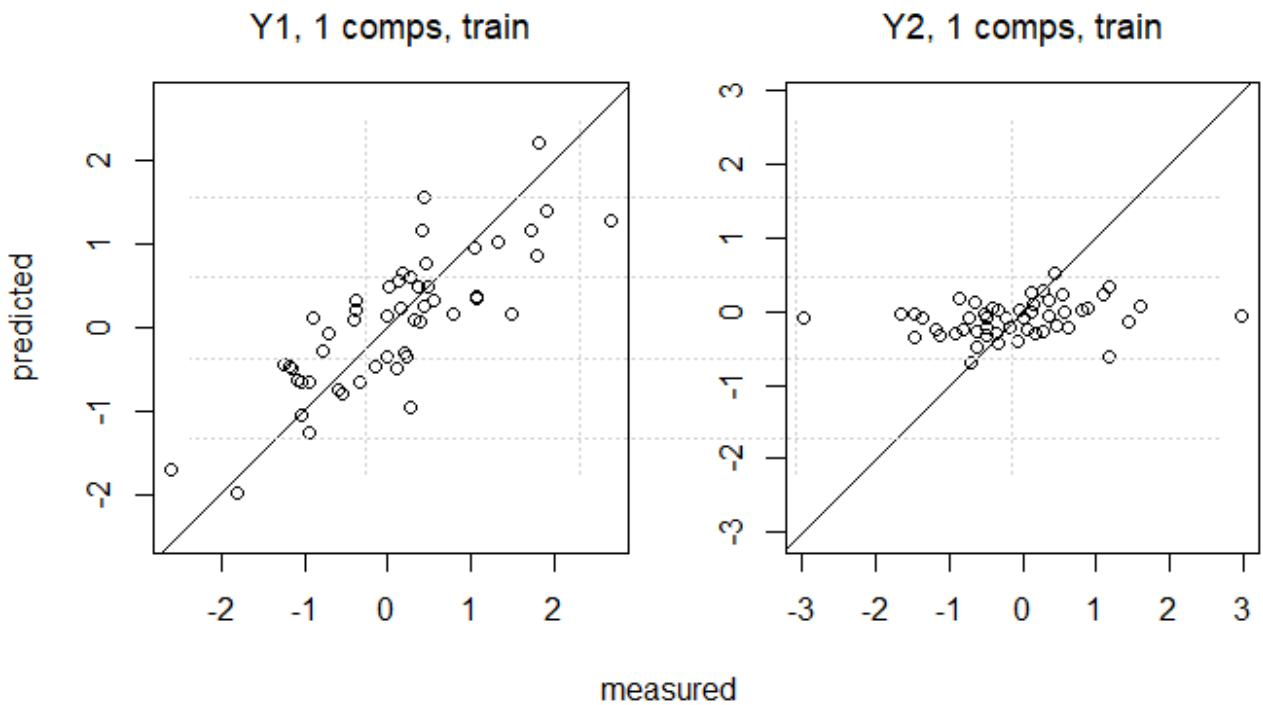
```
1 plot(RMSEP(fit02))
```



フィットの様子

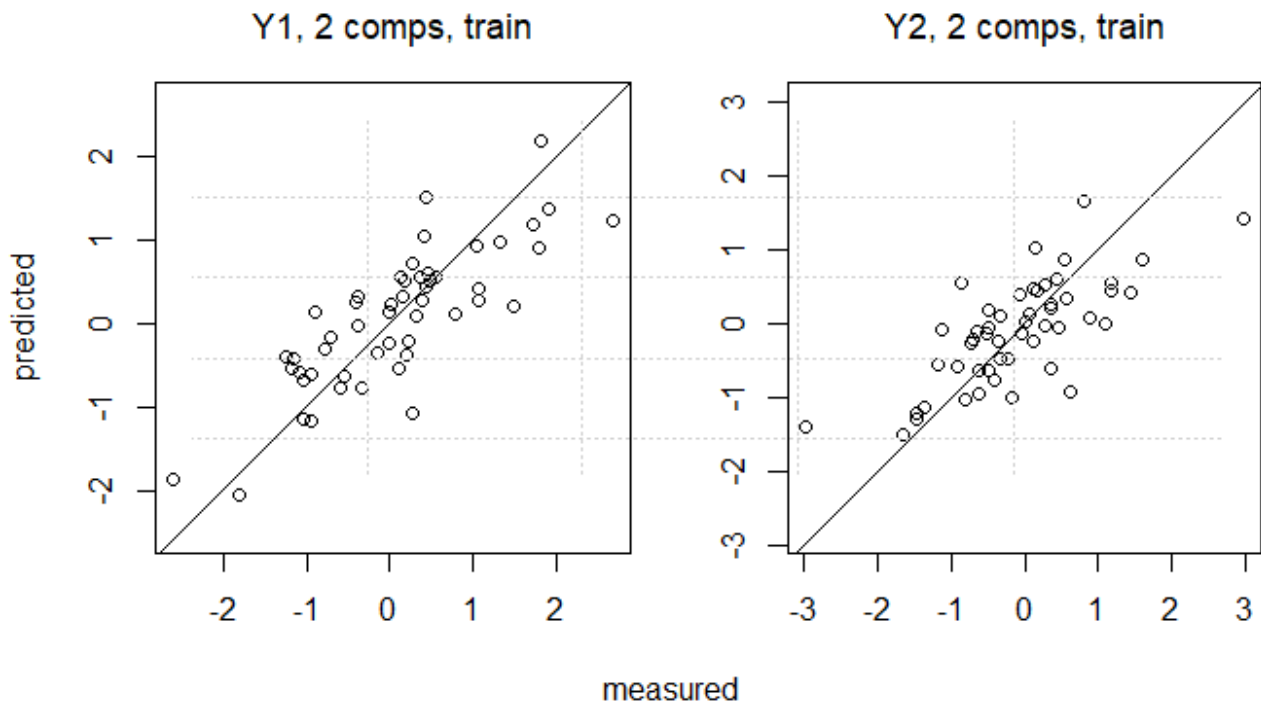
第1主成分だけで予測した場合

```
1 plot(fit02,ncomp=1,line=TRUE,asp=1)
2 grid()
```



第2主成分も使って予測した場合

```
1 plot(fit02,ncomp=2,line=TRUE,asp=1)
2 grid()
```



主成分のスコア

X の主成分のスコアは（最初の6個だけ）

```
1 px2<-scores(fit02)
2 px2%>%
3   round(4)%>%
4   head()
```

```
1 ##   Comp 1   Comp 2
2 ## 1 -0.1106 -0.0836
3 ## 2 -1.3586 -0.1048
4 ## 3  1.1375 -0.0624
5 ## 4 -0.2181 -1.3323
6 ## 5 -0.0030  1.1652
7 ## 6  1.1188 -0.6712
```

Y の主成分のスコアは（最初の6個だけ）

```
1 py2<-Yscores(fit02)
2 py2%>%
3   round(4)%>%
4   head()
```

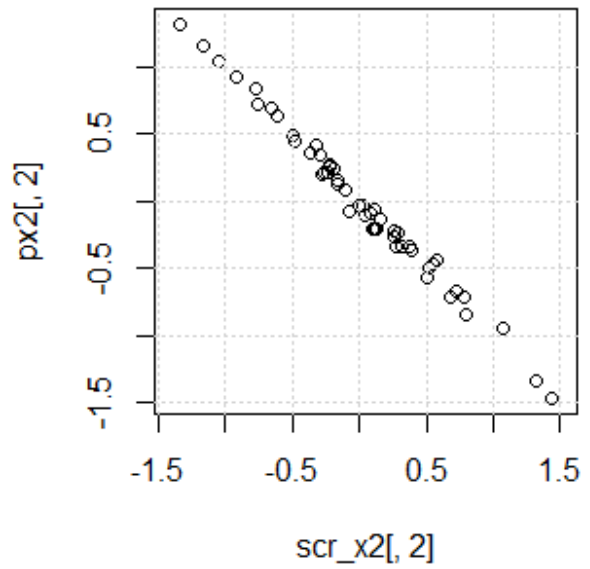
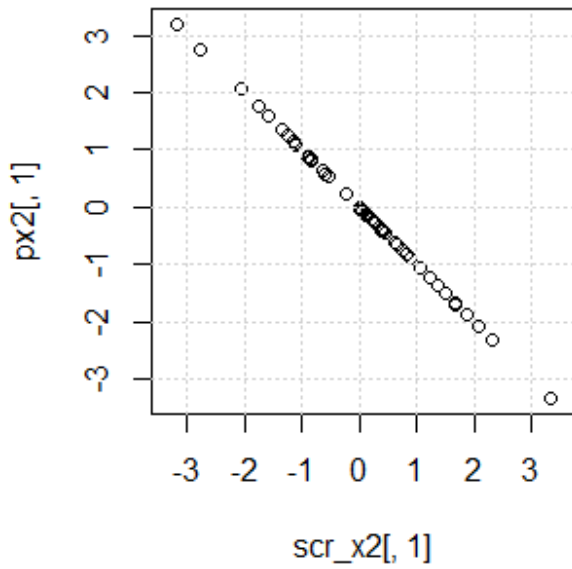
```
1 ##   Comp 1   Comp 2
2 ## 1  0.0619 -0.0882
3 ## 2 -1.5840 -0.1580
4 ## 3  1.7078 -0.0184
5 ## 4 -0.6398 -2.7412
6 ## 5  0.7637  2.5647
7 ## 6  0.4764 -0.3743
```

X のスコアは Xu^* から出した `scr_x2` に（符号を除いて）ほぼ一致する

```

1 par(mfrow=c(1,2))
2 plot(px2[,1]~scr_x2[,1],asp=1)
3 grid()
4 plot(px2[,2]~scr_x2[,2],asp=1)
5 grid()

```

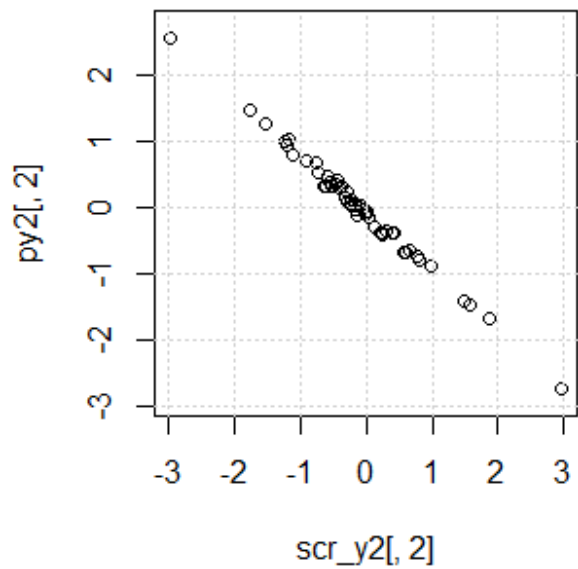
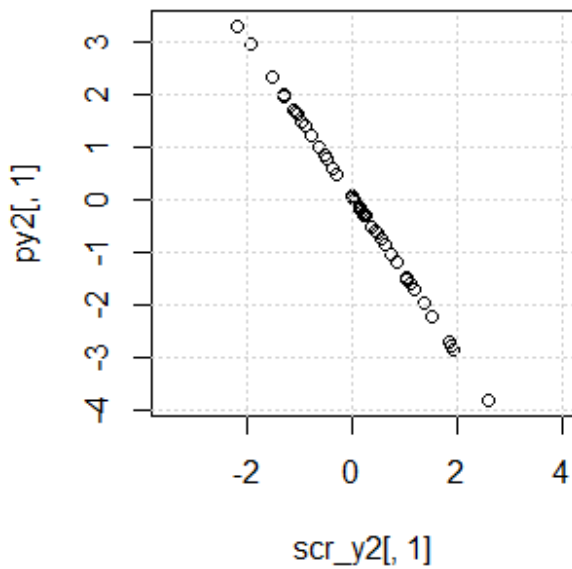


しかし、Yのスコアは Yv^* とスケールが一致しない！

```

1 par(mfrow=c(1,2))
2 plot(py2[,1]~scr_y2[,1],asp=1)
3 grid()
4 plot(py2[,2]~scr_y2[,2],asp=1)
5 grid()

```



やはりYのスコアの分散の調整が行われている。

```
1 lm(Yscores(fit02)~scores(fit02)-1)
```

```
1 ##
2 ## Call:
3 ## lm(formula = Yscores(fit02) ~ scores(fit02) - 1)
4 ##
5 ## Coefficients:
6 ##              Comp 1      Comp 2
7 ## scores(fit02)Comp 1  1.000e+00 -8.270e-17
8 ## scores(fit02)Comp 2  2.075e-01  1.000e+00
```

Yのスコアの分散は

```
1 py2_v<-var(py2)
2 py2_v%>%
3   round(4)
```

```
1 ##      Comp 1 Comp 2
2 ## Comp 1 2.4023 0.1101
3 ## Comp 2 0.1101 0.7093
```

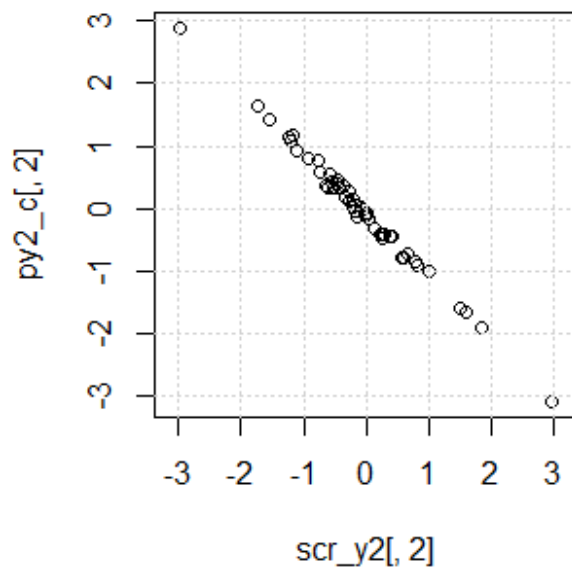
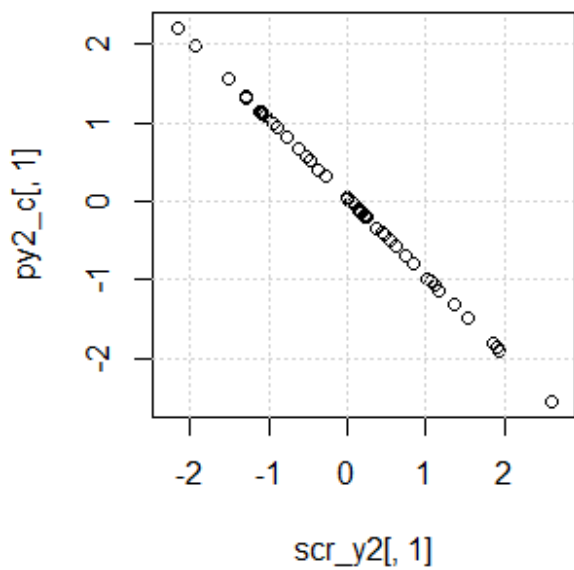
X'Yの最大化から求めたYのスコアの分散は

```
1 xy2_v<-var(scr_y2)
2 xy2_v%>%
3   round(4)
```

```
1 ##      [,1] [,2]
2 ## [1,] 1.0801 0.0414
3 ## [2,] 0.0414 0.9009
```

`pls()` の2つのスコアを、それぞれ割ってやると、X'Yの最大化から求めたYのスコアと一致する。

```
1 py2_c<-sweep(py2,2,
2   c(sqrt(py2_v[1,1]/xy2_v[1,1]),sqrt(py2_v[2,2]/xy2_v[2,2])),
3   FUN="/")%>%
4   round(4)
5 par(mfrow=c(1,2))
6 plot(py2_c[,1]~scr_y2[,1],asp=1)
7 grid()
8 plot(py2_c[,2]~scr_y2[,2],asp=1)
9 grid()
```

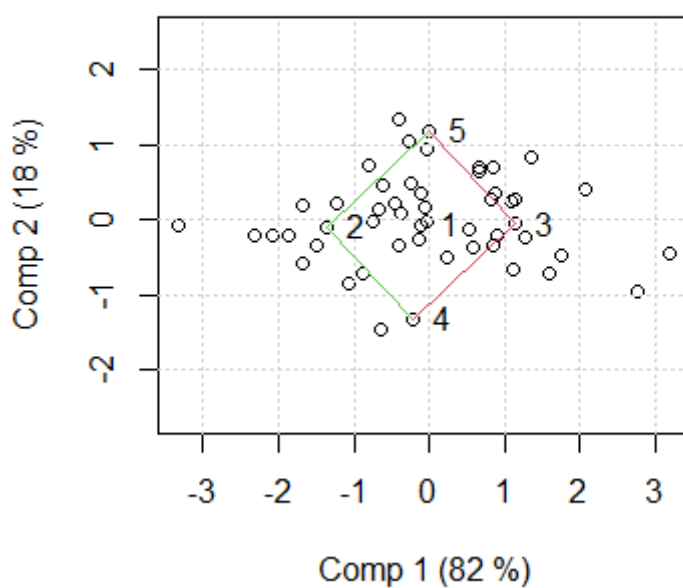
Xのスコアのプロット

`scoreplot()` でもOK

```

1 plot(fit02,plottype="scores",comps=1:2,asp=1)
2 grid()
3 text(x=px2[1:5,1],y=px2[1:5,2],1:5,pos=4)
4 f_grid1(px2)

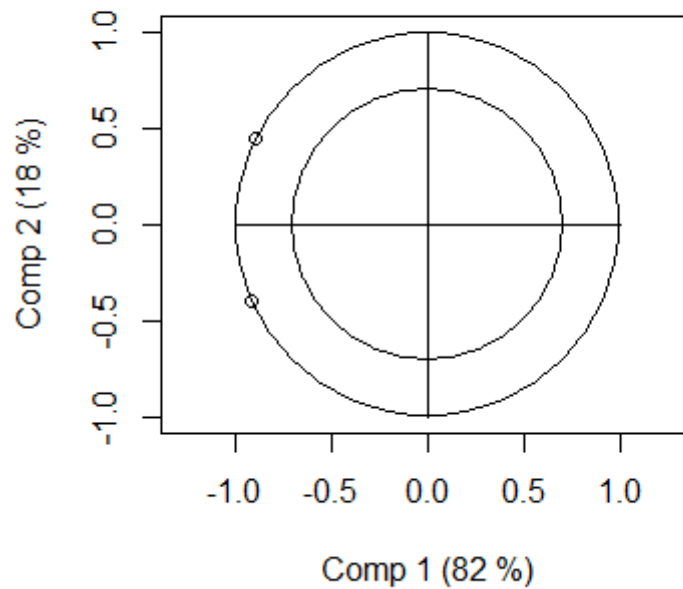
```



スコアとXとの相関のプロット

`plot()` 関数の `plotype="correlation"` 引数でも同じ.

```
1 corrplot(fit02)
```



```
1 cor(x2,px2)
```

```
1 ##          Comp 1    Comp 2
2 ## [1,] -0.9180898 -0.3963724
3 ## [2,] -0.8958415  0.4443737
```

・・・というのをプロットしているようだ. 左下が x_1 で, 左上が x_2 .

x_1 も x_2 も, 第1主成分の方に相関が高そう・・・ということはわかる.

u の出力は `loadings()` で.

```
1 loadings(fit02)
```

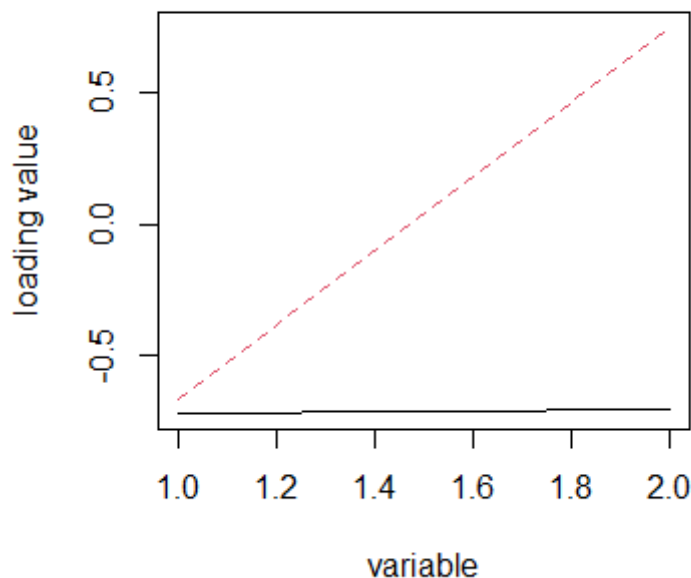
```
1 ##
2 ## Loadings:
3 ##   Comp 1 Comp 2
4 ## x1 -0.716 -0.666
5 ## x2 -0.699  0.746
6 ##
7 ##           Comp 1 Comp 2
8 ## SS loadings  1.002  1.000
9 ## Proportion Var 0.501  0.500
10 ## Cumulative Var 0.501  1.001
```

Comp 1 は, $-0.716*x_1+0.699*x_2$ だよと言っている.

Comp 2 は, $-0.666*x_1+0.746*x_2$ だよと言っている.

2変数だけなので, 図示する意味はあまりないが, `loadingplot()` でいける

```
1 loadingplot(fit02)
```



x_1, x_2 の1, 2を数字として扱っているため, 途中に `1.2` とか余計な目盛が入っているが, これは無視する.

v は, `Yloadings()` で得られるはずだが

```
1 Yloadings(fit02)
```

```
1 ##
2 ## Loadings:
3 ##   Comp 1 Comp 2
4 ## Y1 -0.644  0.180
5 ## Y2 -0.188 -1.114
6 ##
7 ##               Comp 1 Comp 2
8 ## SS loadings   0.450  1.272
9 ## Proportion Var 0.225  0.636
10 ## Cumulative Var 0.225  0.861
```

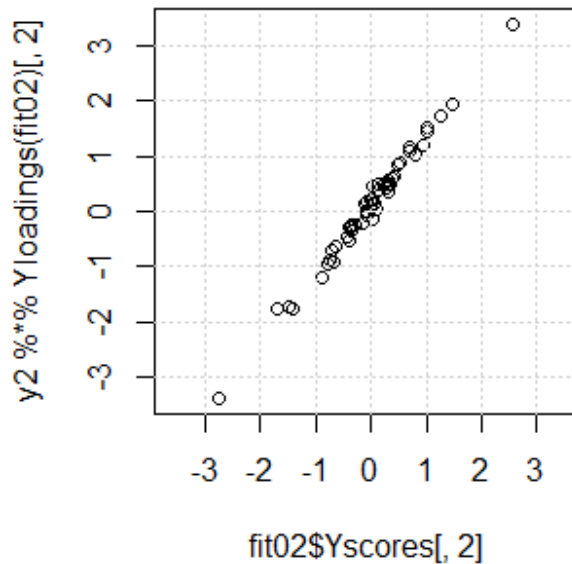
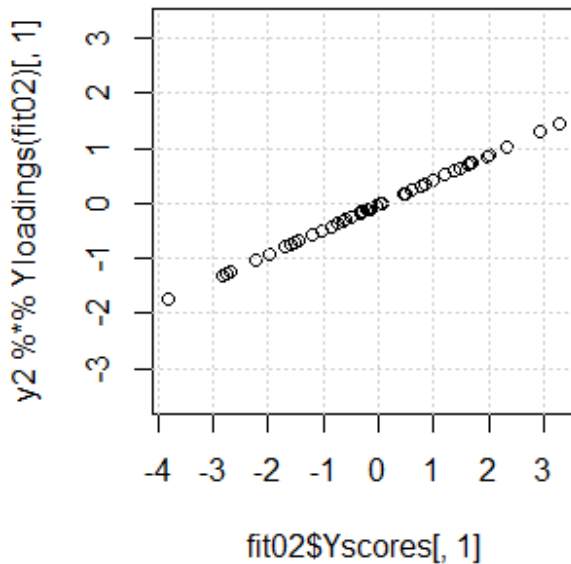
・・・と, loadingが1を超えていたりとよくわからない???

Y との積をとれば, ちゃんとスコアと比例はするが, スケールが全然違う.

```

1 par(mfrow=c(1,2))
2 plot(y2%*%Yloadings(fit02)[,1]~fit02$Yscores[,1],asp=1)
3 grid()
4 plot(y2%*%Yloadings(fit02)[,2]~fit02$Yscores[,2],asp=1)
5 grid()

```



$X'Y$ による Y のスコア Yv^* の分散は

```

1 svd_xy2<-t(x2)%*%y2%>%
2   svd()
3 v_scr_xy2<-y2%*%svd_xy2$v%>%
4   var()
5 v_scr_xy2

```

```

1 ##           [,1]      [,2]
2 ## [1,] 1.08010171 0.04140877
3 ## [2,] 0.04140877 0.90089439

```

`pls()` で出したスコアの分散は、 Y が調整されている分だけ変わっている。

```

1 v_scr_pls2<-fit02$Yscores%>%
2   var()
3 v_scr_pls2

```

```

1 ##           Comp 1    Comp 2
2 ## Comp 1 2.4023141 0.1100564
3 ## Comp 2 0.1100564 0.7093003

```

PLSのLoadingsに、第1主成分の分散の比の平方根を掛けてやると、 v^* の1列目に一致する。

```

1 sqrt(v_scr_pls2[1,1]/v_scr_xy2[1,1])*Yloadings(fit02)[1,1]

```

```
1 ## [1] -0.9599648
```

```
1 sqrt(v_scr_pls2[1,1]/v_scr_xy2[1,1])*Yloadings(fit02)[2,1]
```

```
1 ## [1] -0.2801207
```

しかし、第2主成分の分散の比の平方根を掛けても、 v^* の2列目と近くはなるが、全くは一致しない。

```
1 sqrt(v_scr_pls2[2,2]/v_scr_xy2[2,2])*Yloadings(fit02)[1,2]
```

```
1 ## [1] 0.159741
```

```
1 sqrt(v_scr_pls2[2,2]/v_scr_xy2[2,2])*Yloadings(fit02)[2,2]
```

```
1 ## [1] -0.9880762
```

```
1 v2
```

```
1 ##          [,1]      [,2]
2 ## [1,] 0.9599648 -0.2801207
3 ## [2,] 0.2801207  0.9599648
```

無相関の説明変数を加える

以上、 2×2 変数の場合で説明したが、PLSの強みは、数多くの説明変数の中から相関する主成分を抽出してこること。

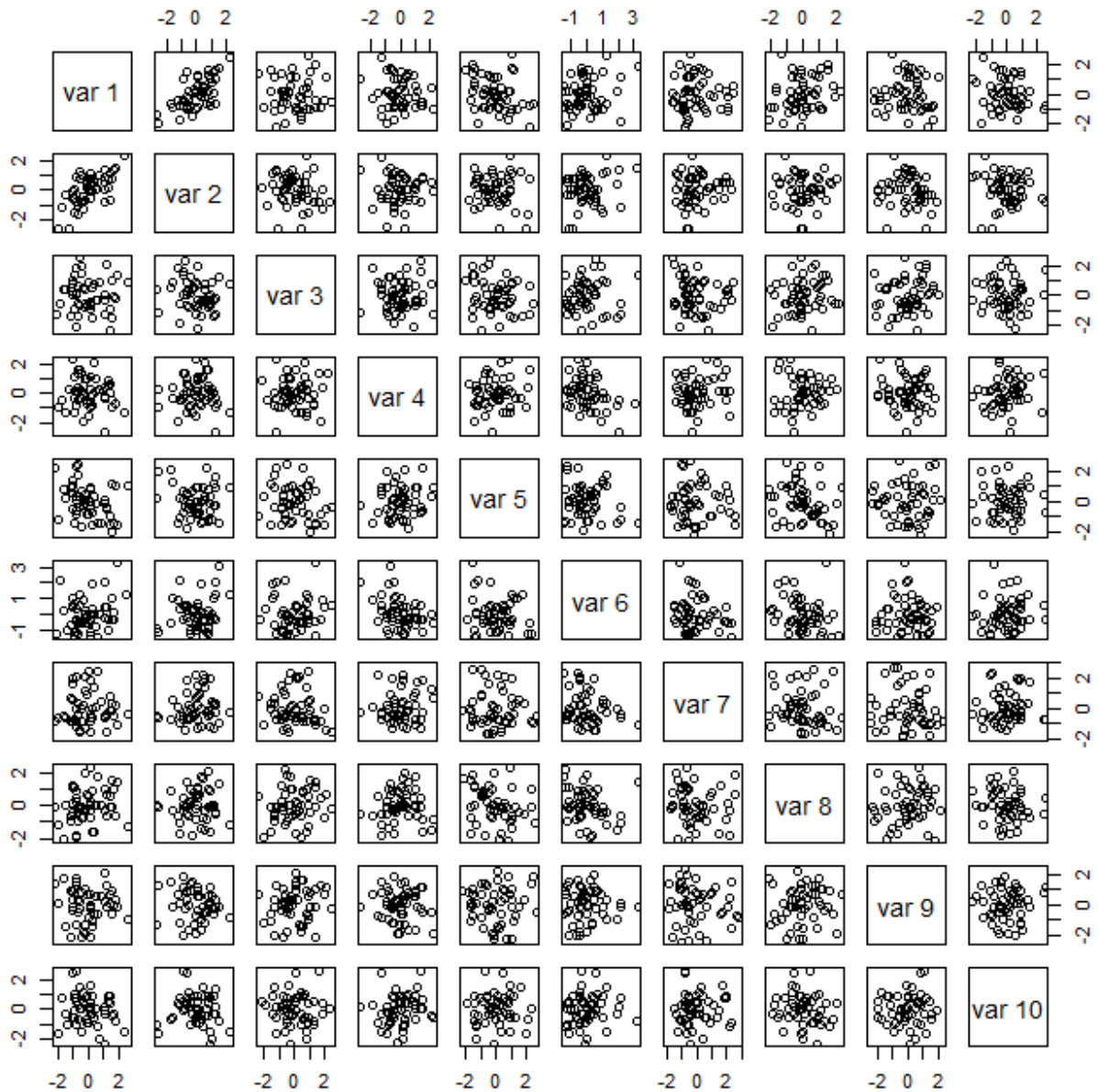
8つのランダムな変数を X に加えてみる。

データ作成

```
1 set.seed(123)
2 x3_c<-rnorm(51*8)%>%
3   matrix(ncol=8)%>%
4   scale()
5 x3<-cbind(x2,x3_c)
```

相関を確認してみる

```
1 pairs(x3,asp=1)
```



x_1, x_2 だけが相関していることを確認.

Y はそのまま使う

```
1 y3<-y2
```

`plsr()` 関数用のデータセット

```
1 d03<-data.frame(x=I(x3),y=I(y3))
```

`plsr()` 関数で

```
1 fit03<-plsr(y~x,
2           ncomp=10,
3           data=d03,
4           validation="CV")
```

要約

```
1 summary(fit03)
```

```
1  ## Data:   X dimension: 51 10
2  ## Y dimension: 51 2
3  ## Fit method: kernelpls
4  ## Number of components considered: 10
5  ##
6  ## VALIDATION: RMSEP
7  ## Cross-validated using 10 random segments.
8  ##
9  ## Response: Y1
10 ##      (Intercept) 1 comps 2 comps 3 comps 4 comps 5 comps 6 comps
11 ## CV      1.032  0.7149  0.6844  0.6857  0.7112  0.6963  0.7012
12 ## adjCV    1.032  0.7059  0.6780  0.6798  0.6996  0.6886  0.6935
13 ##      7 comps 8 comps 9 comps 10 comps
14 ## CV      0.7109  0.7188  0.7215  0.7220
15 ## adjCV    0.7029  0.7100  0.7125  0.7129
16 ##
17 ## Response: Y2
18 ##      (Intercept) 1 comps 2 comps 3 comps 4 comps 5 comps 6 comps
19 ## CV      0.9777  0.9649  0.8875  0.8269  0.8100  0.7647  0.7612
20 ## adjCV    0.9777  0.9642  0.8801  0.8096  0.8024  0.7559  0.7520
21 ##      7 comps 8 comps 9 comps 10 comps
22 ## CV      0.7596  0.7603  0.7599  0.7594
23 ## adjCV    0.7506  0.7513  0.7509  0.7506
24 ##
25 ## TRAINING: % variance explained
26 ##      1 comps 2 comps 3 comps 4 comps 5 comps 6 comps 7 comps 8 comps
27 ## X      19.526  34.28  40.26  50.45  60.00  70.42  76.90  85.01
28 ## Y1     60.092  66.74  67.06  69.49  70.04  70.06  70.07  70.09
29 ## Y2      9.682  38.86  58.05  59.00  61.19  61.63  61.69  61.69
30 ##      9 comps 10 comps
31 ## X      94.23  100.00
32 ## Y1     70.09  70.09
33 ## Y2     61.69  61.69
```

主成分数をいくつにするか決める

その前にRMSEPって？

Yの実際値と推定値との誤差の平方平均の平方根がRMSEP。
しかし、計算してみると一致しない。

第1主成分までのRMSEPは

```
1 RMSEP(fit03,ncomp=1)
```

```

1 ##
2 ### Response: Y1
3 ##      (Intercept)  1 comps
4 ## CV          1.032  0.7149
5 ## adjCV       1.032  0.7059
6 ##
7 ### Response: Y2
8 ##      (Intercept)  1 comps
9 ## CV          0.9777 0.9649
10 ## adjCV       0.9777 0.9642

```

ところが、残差の平方平均の平方根は

```

1 fit03$residuals[,1]^2>%
2   apply(.,2,mean)>%
3   sqrt()

```

```

1 ##      Y1      Y2
2 ## 0.6390451 0.9109793

```

なんかちょっと違う。

実際には、クロスバリデーションで得られたRMSEPを計算している。

すべてのサンプルを使って推定すると（引数 `estimate="all"` で、`train` の `1 comp` のところを見ると・・・），単純に計算したRMSEPと一致する。

```

1 RMSEP(fit03,ncomp=1,estimate="all")

```

```

1 ##
2 ### Response: Y1
3 ##      (Intercept)  1 comps
4 ## train          1.012  0.6390
5 ## CV            1.032  0.7149
6 ## adjCV        1.032  0.7059
7 ##
8 ### Response: Y2
9 ##      (Intercept)  1 comps
10 ## train          0.9586 0.9110
11 ## CV            0.9777 0.9649
12 ## adjCV        0.9777 0.9642

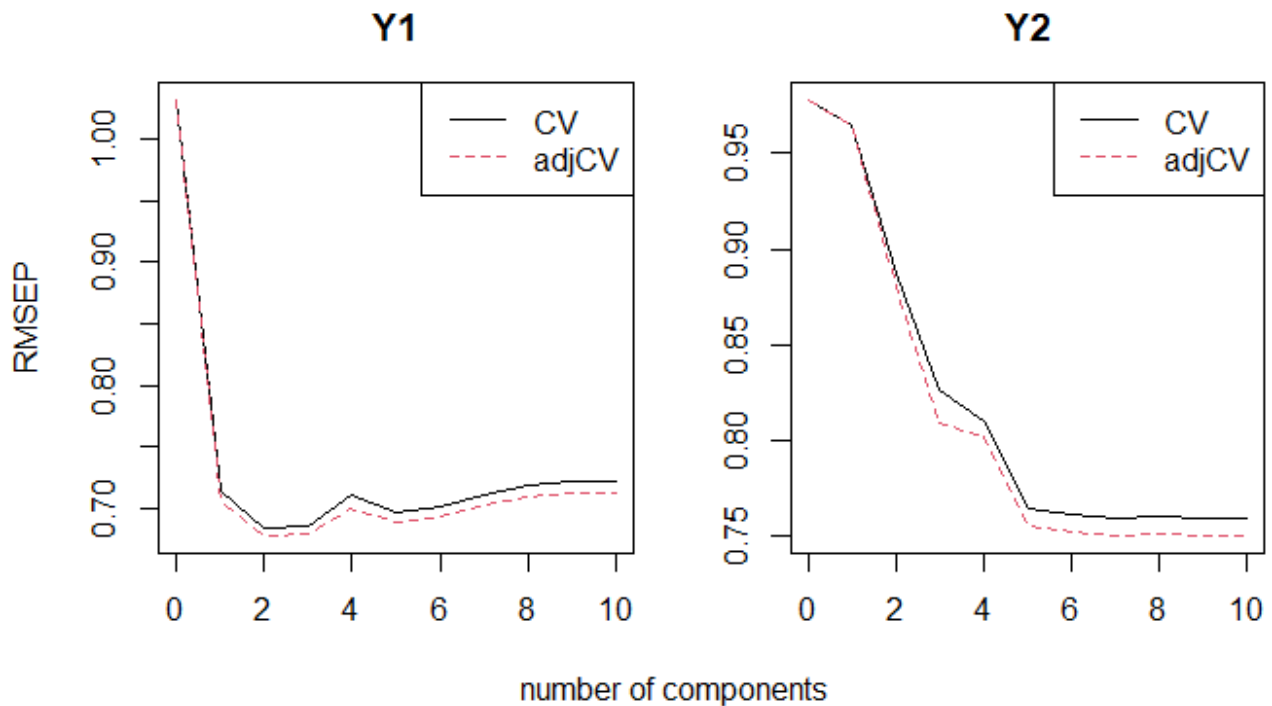
```

RMSEPで

```

1 plot(RMSEP(fit03), legendpos = "topright")

```

成分数を選ぶヒューリスティクスとして、`selectNcomp()` 関数が用意されている。使てみたいが、これはYが1変数のときしか使えないようだ。

どんなものか知るだけのために、 y_1 だけにPLS回帰してみよう。

```
1 fit03_1<-plsr(y[,1]~x,ncomp=10,data=d03,validation="CV")
```

`selectNcomp()` 関数を使ってみる。判断基準は2通りあり、引数 `method="onsigma"` と引数 `method="randomization"` で使い分ける。

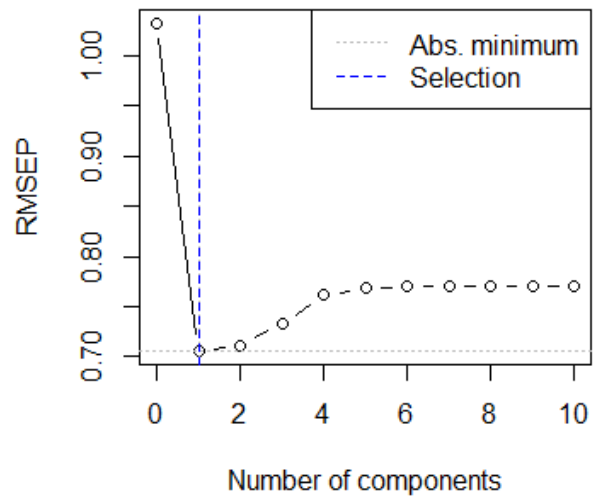
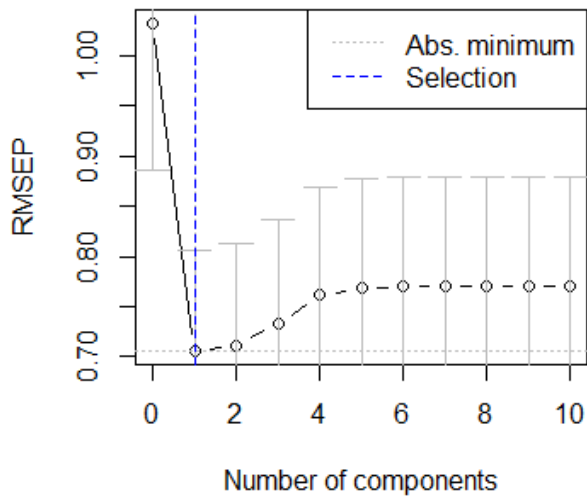
"onsigma" は、まず、RMSEPが最低の点を求め、そこから1標準誤差以内で最も少ない主成分数を採用するという方法。

"randomization" は、やはりRMSEPが最低の点を求め、そこから有意差がない主成分数を採用するという方法。

```
1 par(mfrow=c(1,2))
2 selectNcomp(fit03_1, method = "onsigma",plot=TRUE)
```

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

```
1 selectNcomp(fit03_1, method = "randomization",plot=TRUE)
```



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

両方とも1主成分が最適と出してきた。

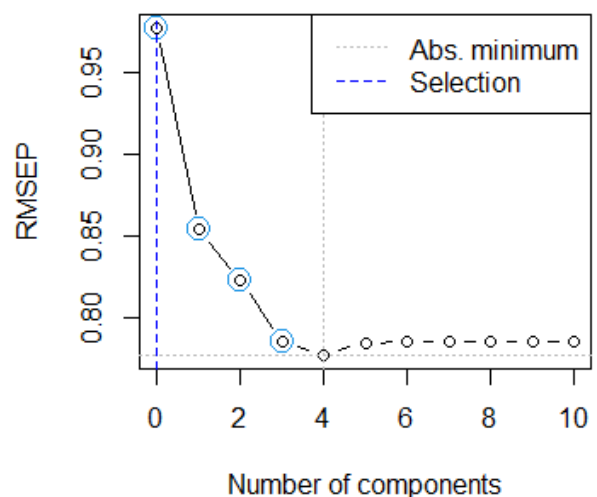
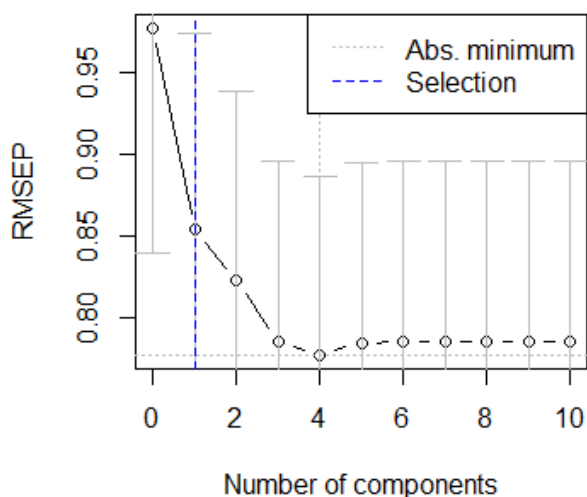
y_2 についてもやってみよう。

```
1 fit03_2<-plsr(y[,2]~x,ncomp=10,data=d03,validation="CV")
```

```
1 par(mfrow=c(1,2))
2 selectNcomp(fit03_2, method = "onesigma",plot=TRUE)
```

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

```
1 selectNcomp(fit03_2, method = "randomization",plot=TRUE)
```



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

1か0なのだそう。

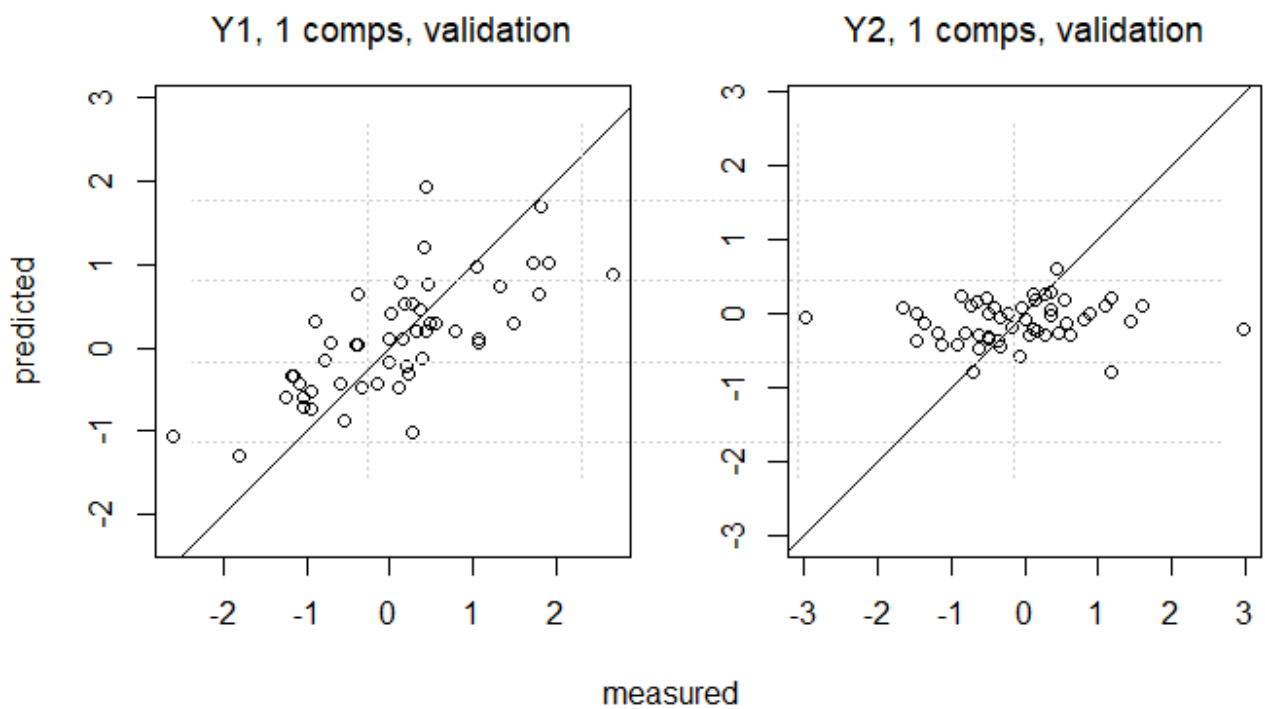
しかし、これはあくまで、 Y を直接対応させた場合の結果。

予測値で

もとの Y と予測値とをプロットして確認

1 主成分だけだと y_1 はいい感じでfitしているが、 y_2 はいまいち。

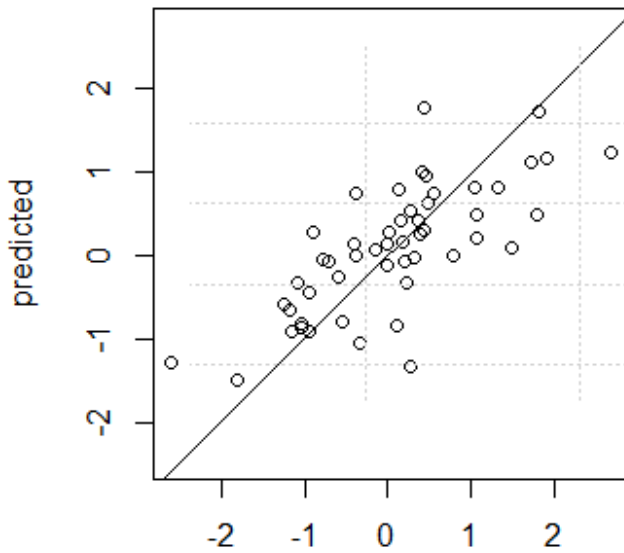
```
1 plot(fit03, ncomp = 1, asp = 1, line = TRUE)
2 grid()
```



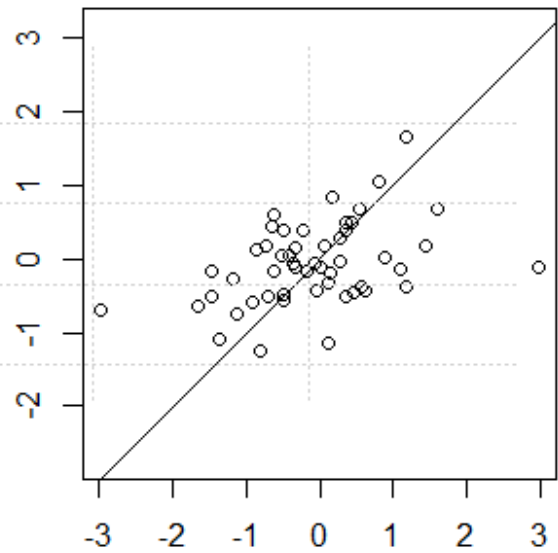
第2主成分まで使うと y_1 のfitはあんまりかわらんが、 y_2 はだいぶよくなった。

```
1 plot(fit03, ncomp = 2, asp = 1, line = TRUE)
2 grid()
```

Y1, 2 comps, validation



Y2, 2 comps, validation

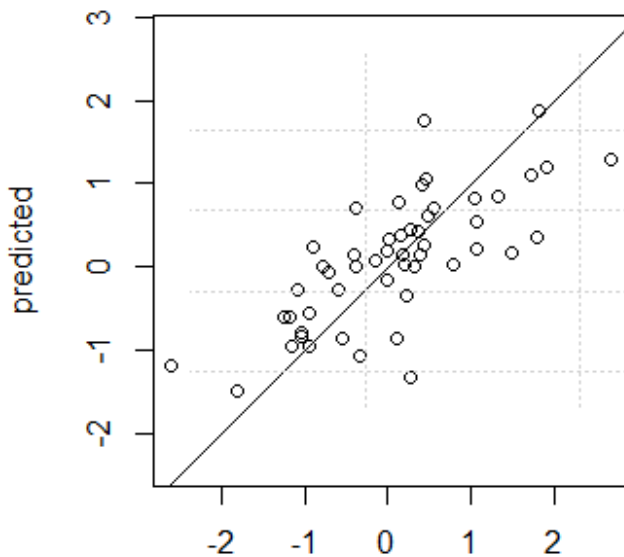


measured

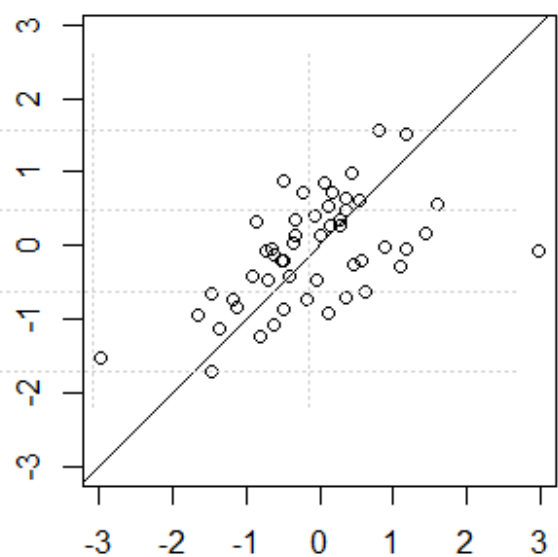
第3主成分まで使うと、両者ともあんまり変わらない。

```
1 plot(fit03, ncomp = 3, asp = 1, line = TRUE)
2 grid()
```

Y1, 3 comps, validation



Y2, 3 comps, validation



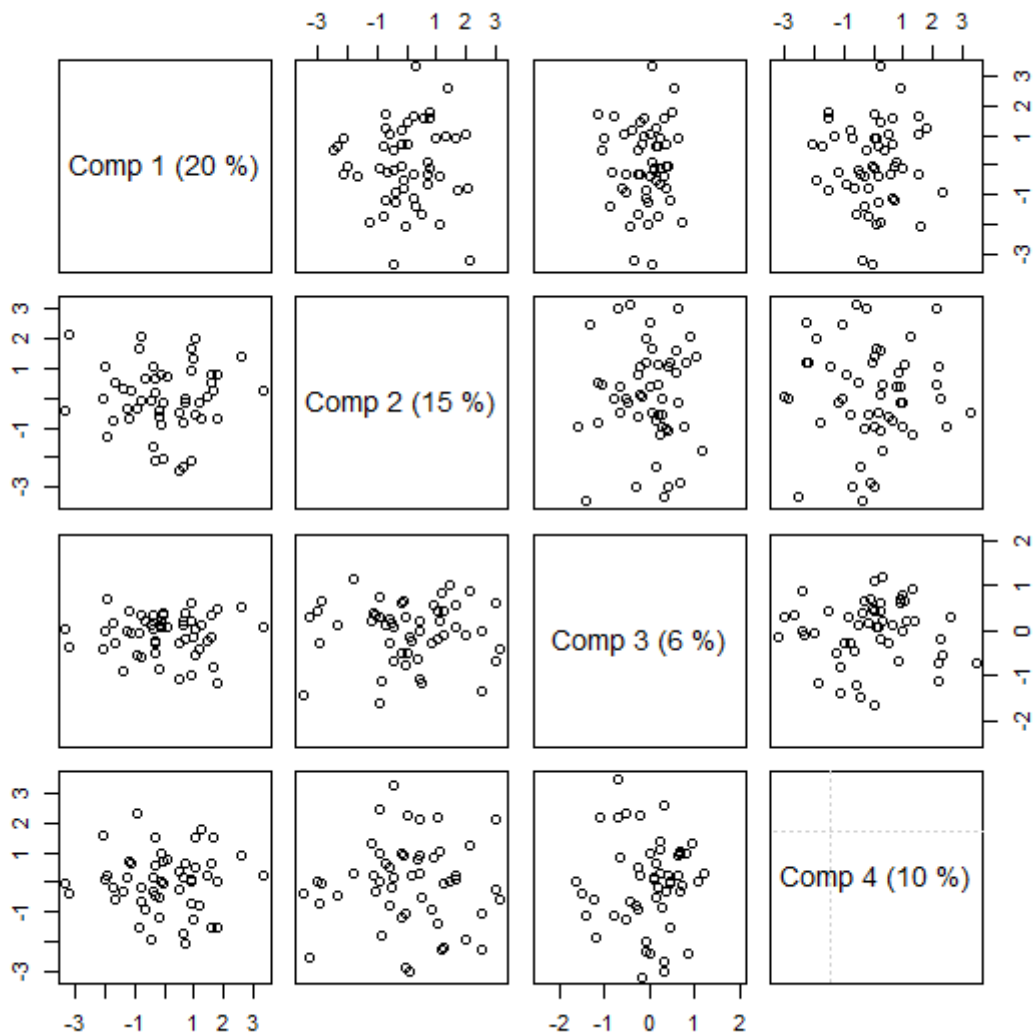
measured

2主成分ということか？

Xのスコアのプロット

Xの第4主成分までスコアのプロットしてみると...

```
1 scoreplot(fit03, asp=1, comps = 1:4)
2 grid()
```



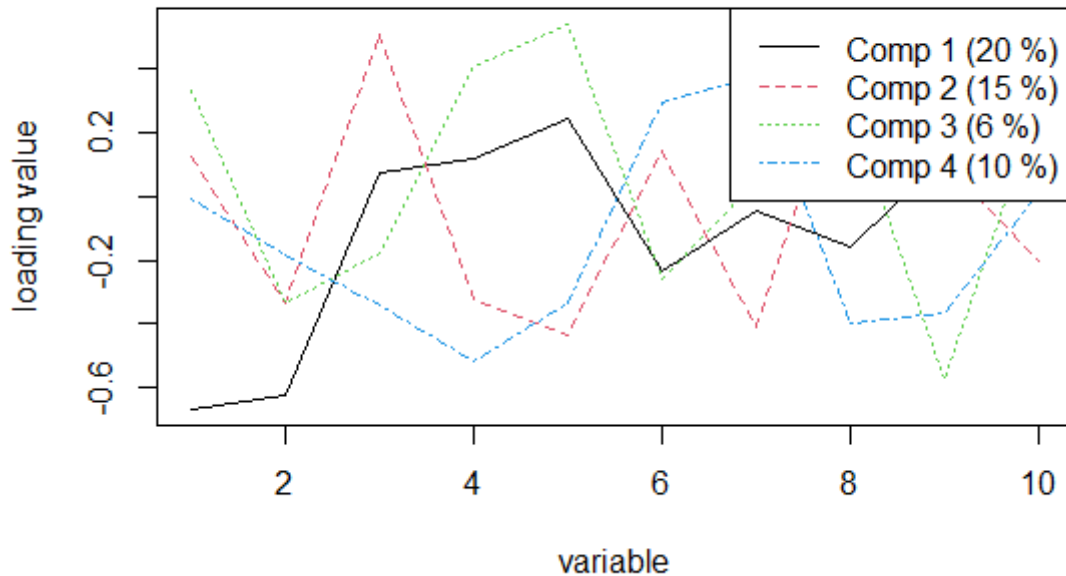
う～ん、この情報はどうなんだろう。

説明力の高い X の変数を選択する

第1主成分だけでいくとして、 X のうち、どの変数を説明変数として採用するか・・・

u (loading value) を確認する

```
1 | loadingplot(fit03, legendpos = "topright", comps=1:4)
```



Comp 1 の loading value を見ると、変数 1 と 2 だけが大きくぶれている。

Comp 2 の loading value は、特に大きな値はないなあ。

これを見る限り、 x_1 と x_2 だけで 1 つの主成分を抽出し、 y_1, y_2 を説明するのが適切か？ という仮説が出される・・・たぶん。

ダメ押しで、スコアと X との相関も。

次は、第 1 主成分と第 2 主成分のスコアと X の各変数との相関

loading

value のノルムで、左端の x_1, x_2 は 1 に近いが、他は 0.5 あたりか、それ未満。あんまり相関していないということ。

補足:シミュレーション

係数の不偏性を確かめる。

51×2 の X, Y について、OLS ($Y \leftarrow X$) と PCR ($Y \leftarrow X$ のスコア)、それに $X'Y$ の特異値分解から求めたスコア間の OLS 回帰の結果を比較する。

3 通りの係数を計算する関数

とりあえず 1 回の推定をする関数を定義する。

攪乱項は X と Y の両方に与え、その標準偏差は両者で同じとする。
各推定 4 つずつ、3 種類で、12 個の推定係数が出力。

```

1 f_pls01<-function(x1,y1,sd1=0.1){
2   x2<-x1+matrix(rnorm(nrow(x1)*2,mean=0,sd=sd1),ncol=2)

```

```

3 y2<-y1+matrix(rnorm(nrow(y1)*2,mean=0,sd=sd1),ncol=2)
4 fit_OLS<-lm(y2~x2-1)%>%
5   coef()
6 prc01<-prcomp(x1)
7 Px<-prc01$x
8 fit_PCR<-lm(y2~Px-1)%>%
9   coef()
10 svd_xy<-svd(t(x2)%*%y2)
11 Scrx<-x2%*%svd_xy$u
12 Scry<-y2%*%svd_xy$v
13 fit_PLS<-lm(Scry~Scrx-1)%>%
14   coef()
15 coef01<-c(fit_OLS,fit_PCR,fit_PLS)
16 coef01
17 }

```

シミュレーション

攪乱項の標準偏差を0.2として、1000回推定

```

1 fit_sim01<-matrix(nrow=1000,ncol=12)
2 for(i in 1:1000){
3   fit_sim01[i,]<-f_pls01(x1,y1,sd1=0.2)
4 }
5
6 fit_sim_OLS<-fit_sim01[,1:4]
7 fit_sim_PCR<-fit_sim01[,5:8]
8 fit_sim_PLS<-fit_sim01[,9:12]

```

OLS

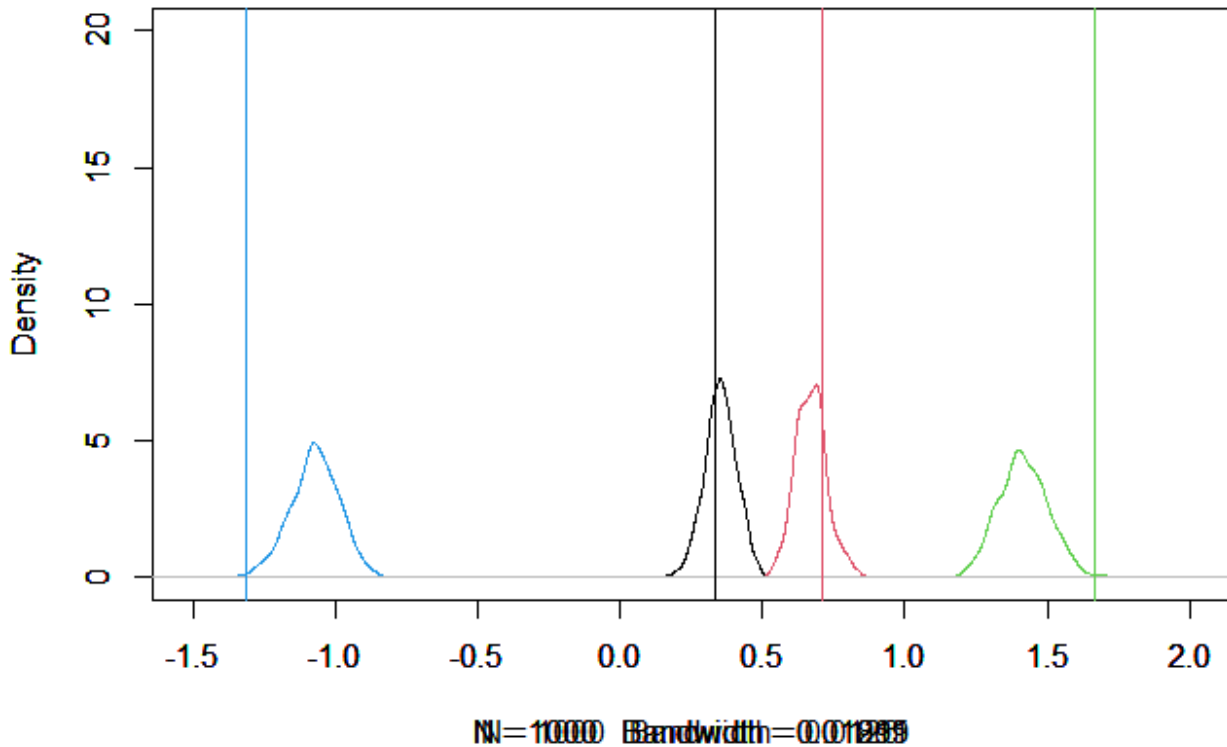
$Y \leftarrow X$ のOLS推定

```

1 fit_OLS0<-a1%*%b0%*%solve(a1)
2 plot(density(fit_sim_OLS[,1]),xlim=c(-1.5,2),ylim=c(0,20),col=1)
3 par(new=T)
4 plot(density(fit_sim_OLS[,2]),xlim=c(-1.5,2),ylim=c(0,20),col=2)
5 par(new=T)
6 plot(density(fit_sim_OLS[,3]),xlim=c(-1.5,2),ylim=c(0,20),col=3)
7 par(new=T)
8 plot(density(fit_sim_OLS[,4]),xlim=c(-1.5,2),ylim=c(0,20),col=4)
9 abline(v=fit_OLS0[1,1],col=1)
10 abline(v=fit_OLS0[2,1],col=2)
11 abline(v=fit_OLS0[1,2],col=3)
12 abline(v=fit_OLS0[2,2],col=4)

```

density.default(x = fit_sim_OLS[, 2])



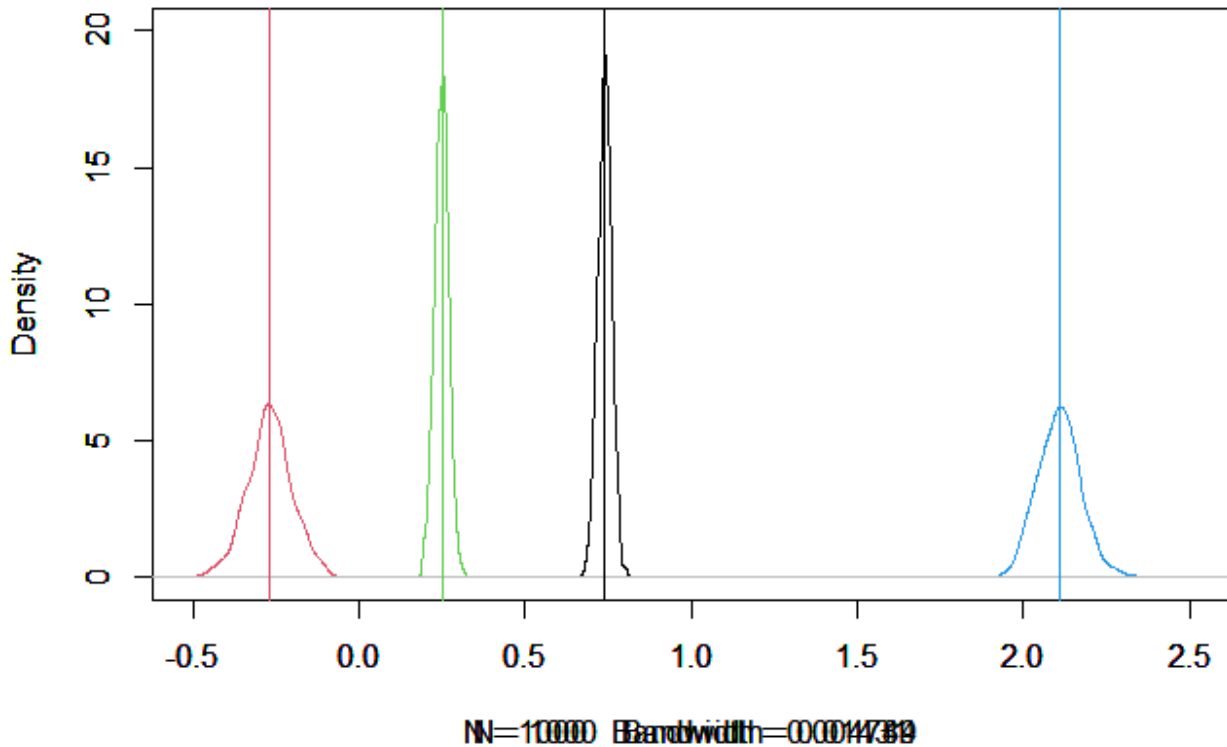
垂直線が真値。不偏性がない！推定値の分散も大きい。

PCR

$Y \leftarrow X$ の主成分分析のスコア

```
1 fit_PCR0<-b0%*%solve(a1)
2 plot(density(fit_sim_PCR[,1]),xlim=c(-0.5,2.5),ylim=c(0,20),col=1)
3 par(new=T)
4 plot(density(fit_sim_PCR[,2]),xlim=c(-0.5,2.5),ylim=c(0,20),col=2)
5 par(new=T)
6 plot(density(fit_sim_PCR[,3]),xlim=c(-0.5,2.5),ylim=c(0,20),col=3)
7 par(new=T)
8 plot(density(fit_sim_PCR[,4]),xlim=c(-0.5,2.5),ylim=c(0,20),col=4)
9 abline(v=fit_PCR0[1,1],col=1)
10 abline(v=fit_PCR0[2,1],col=2)
11 abline(v=fit_PCR0[1,2],col=3)
12 abline(v=fit_PCR0[2,2],col=4)
```


density.default(x = fit_sim_PCR[, 2])



不偏性がありそうだ。

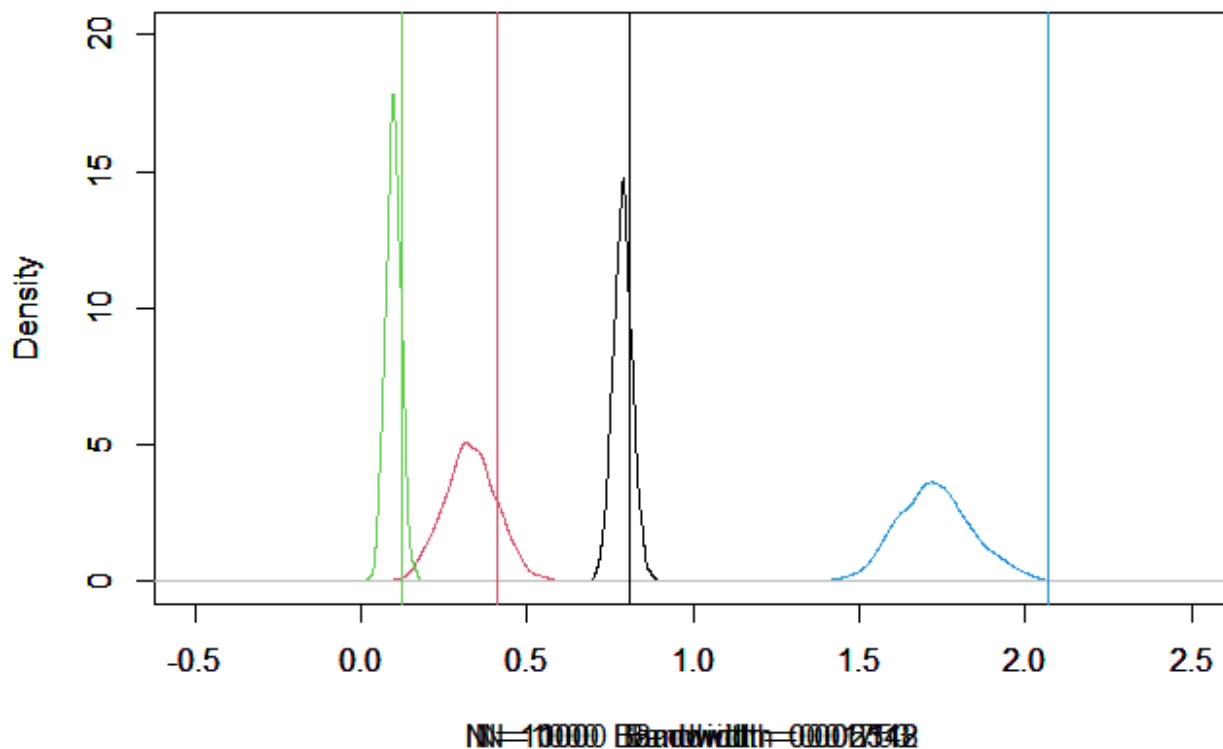
X の第1主成分の係数（黒と緑）は標準誤差が小さいが、第2主成分の標準誤差が長が大きくなっている。

PLS

$X'Y$ の特異値分解から、スコア間の関係 $Yv \leftarrow Xu$

```
1 fit_PLS0 <- xyb1
2 plot(density(fit_sim_PLS[, 1]), xlim=c(-0.5, 2.5), ylim=c(0, 20), col=1)
3 par(new=T)
4 plot(density(fit_sim_PLS[, 2]), xlim=c(-0.5, 2.5), ylim=c(0, 20), col=2)
5 par(new=T)
6 plot(density(fit_sim_PLS[, 3]), xlim=c(-0.5, 2.5), ylim=c(0, 20), col=3)
7 par(new=T)
8 plot(density(fit_sim_PLS[, 4]), xlim=c(-0.5, 2.5), ylim=c(0, 20), col=4)
9 abline(v=fit_PLS0[1, 1], col=1)
10 abline(v=fit_PLS0[2, 1], col=2)
11 abline(v=fit_PLS0[1, 2], col=3)
12 abline(v=fit_PLS0[2, 2], col=4)
```

density.default(x = fit_sim_PLS[, 2])



Xu の第1主成分（黒と緑）については、だいぶいいところまで来ているが、不偏性は怪しい。

第2主成分（赤と青）はだいぶ怪しい???

実際のPLSでは、このように一気に Xu と Yv を算出するのではなく、まず、第1主成分 Xu_1 と Yv_1 を算出し、次いで、それを前提に Xu_2 と Yv_2 を算出・・・というような手順を踏むようだ。（そこから先は、一般的なPLSの解説がたくさんある）

そもそも、PLSの関心が、係数の値そのものにあるのではなく、説明力の高い主成分と、その主成分と相関の高い X の特定にある。`plsr()`関数などは、 X に合わせて Y のスケールそのものを変えてしまう。

複数の Y を主成分として扱えるところに魅力を感じたが、そこはあんまり意味なさそう。

たくさんの X の中から、PLSで主成分とそれに相関の高い X を抽出し、PCRで回帰するのが妥当な使い方のような気がする。

結論

1. PLSは $(Xu)'Yv$ を最大化するように u^* と v^* を決める。
2. この u^* と v^* は、 $X'Y$ の特異値分解から得られる。
3. 推定の考え方がOLSとは違うので、 X と Y にかく乱項が加わった場合、係数の不偏性は怪しい。
4. PLSの実際の利用は、多数の変数を持つ X から Y の説明力の高い X の主成分を抽出し、それと相関の高い X を特定すること。
5. 変数間の係数に関心がある場合は、PLSで特定した X でPCRをするのがよさそう。
6. 実際のPLSは、 $X'Y$ の特異値分解という単純なものでなく、主成分を1つずつ算出していく複雑なアルゴリズムが複数ある。

7. Rのパッケージplsにある `plsr()` 関数で出力される Y のスコア (`Yscores`) は X のスコアに合わせてスケールを調整している所以要注意.