

Rによる判別分析

[本稿のWebページ](#)

古橋 武

2. 線形判別分析(2変量)

2.1 基礎理論

判別分析は判別線(面)を同定する.

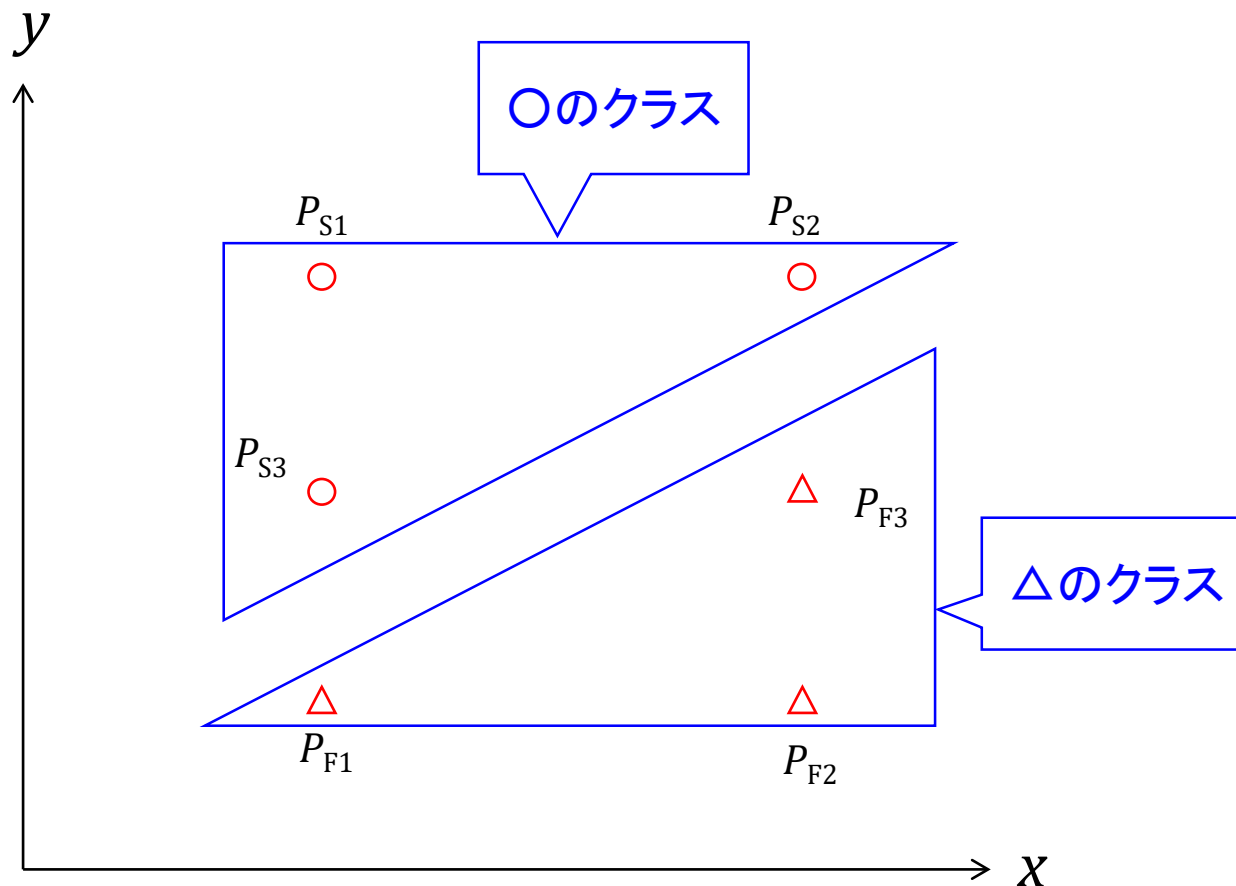


図2.1 判別分析の様子

二つのクラスを分ける線を同定する課題

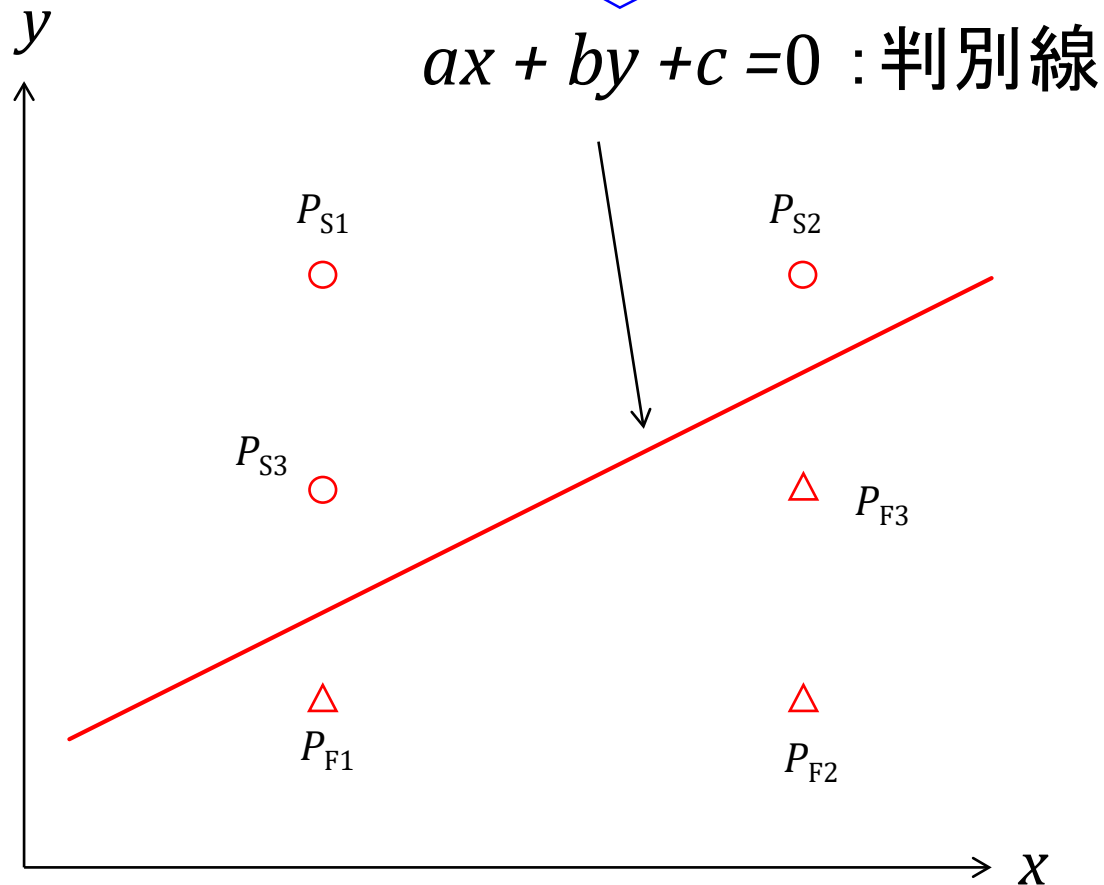
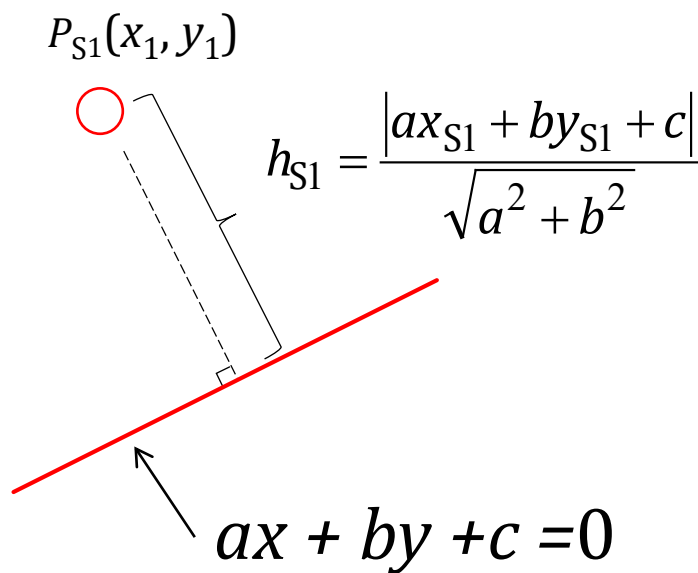


図2.1 判別分析の様子

では、判別線の基準は？

$$h_{S1} = \frac{|ax_{S1} + by_{S1} + c|}{\sqrt{a^2 + b^2}} \quad (2.2)$$



$$a^2 + b^2 = 1$$

とすれば

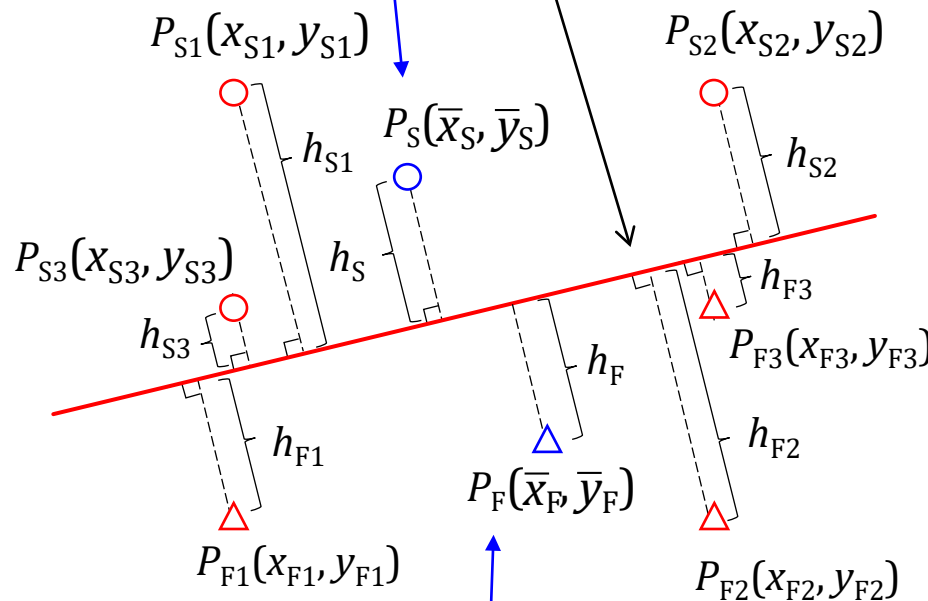
$$h_{S1} = |ax_{S1} + by_{S1} + c|$$

図2.2 データから判別線までの距離

$$\bar{x}_S = \frac{1}{n_S} \sum_{i=1}^{n_S} x_{Si}, \quad \bar{y}_S = \frac{1}{n_S} \sum_{i=1}^{n_S} y_{Si}$$

クラス内散布 s_W^2

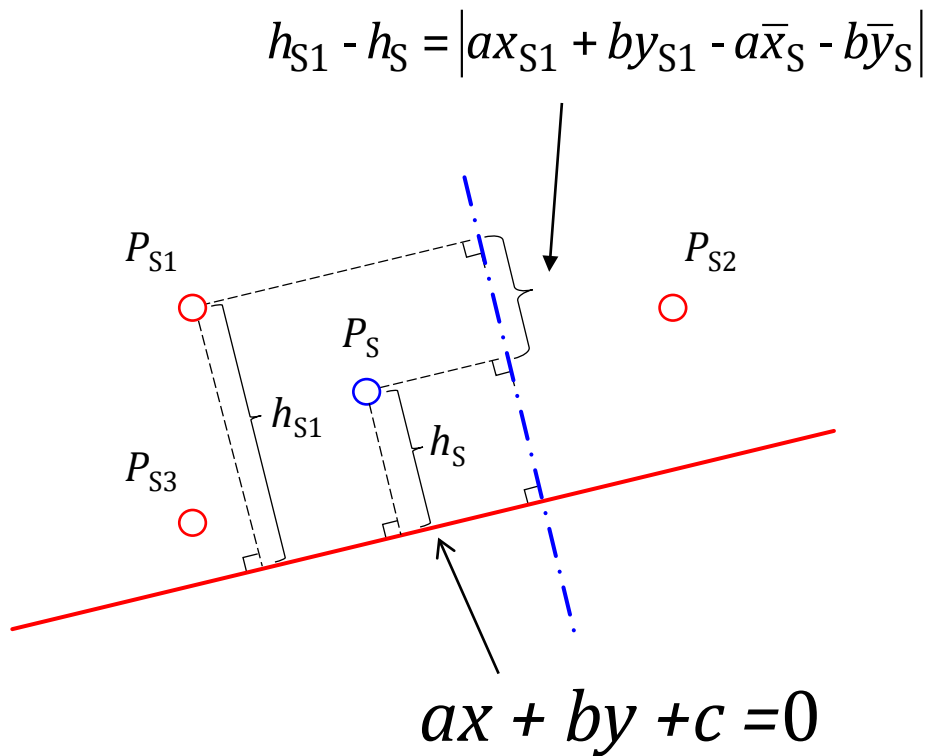
$$ax + by + c = 0$$



$$s_W^2 = (h_{S1} - h_S)^2 + (h_{S2} - h_S)^2 + (h_{S3} - h_S)^2 + (h_{F1} - h_F)^2 + (h_{F2} - h_F)^2 + (h_{F3} - h_F)^2$$

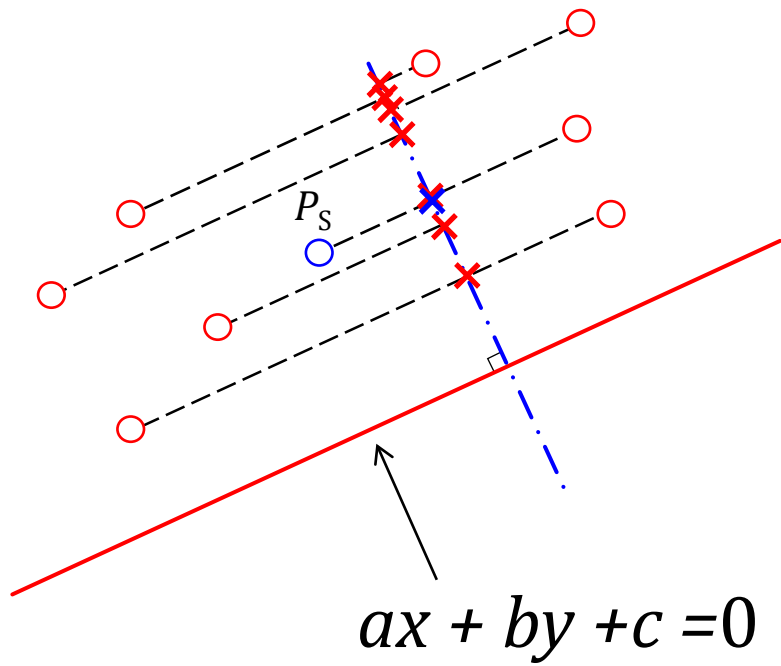
$$\bar{x}_F = \frac{1}{n_F} \sum_{i=1}^{n_F} x_{Fi}, \quad \bar{y}_F = \frac{1}{n_F} \sum_{i=1}^{n_F} y_{Fi}$$

図2.3 各データと判別線までの距離

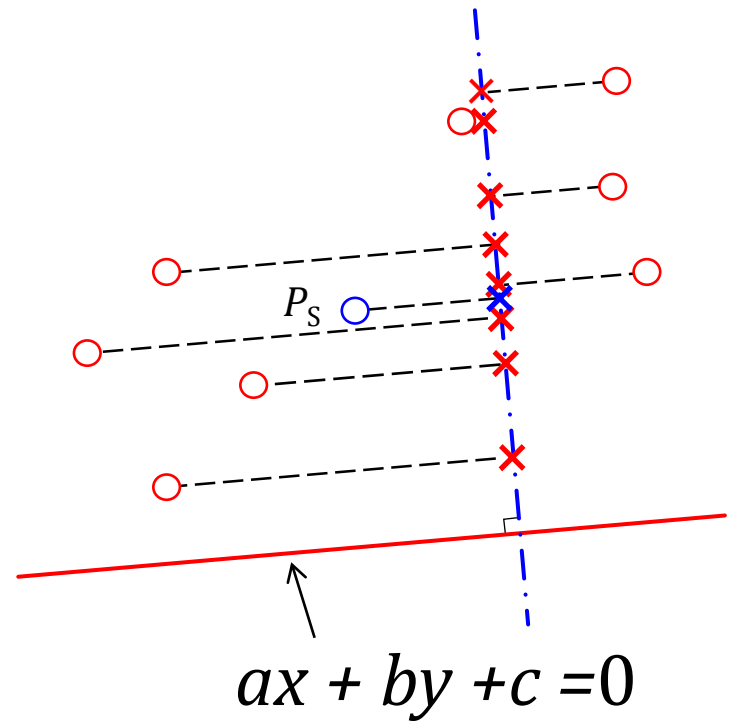


$$\begin{aligned}
 s_W^2 &= (h_{S1} - h_S)^2 \\
 &+ (h_{S2} - h_S)^2 \\
 &+ (h_{S3} - h_S)^2 \\
 &+ (h_{F1} - h_F)^2 \\
 &+ (h_{F2} - h_F)^2 \\
 &+ (h_{F3} - h_F)^2 \\
 &= \sum_{i=1}^{n_S} (ax_{Si} + by_{Si} - a\bar{x}_S - b\bar{y}_S)^2 \\
 &+ \sum_{i=1}^{n_F} (ax_{Fi} + by_{Fi} - a\bar{x}_F - b\bar{y}_F)^2
 \end{aligned}$$

図2.4 ○のクラス内のデータの判別線との距離の差分



(a) クラス内散布小



(b) クラス内散布大

図2.5 判別線の傾きとクラス内散布

全体の平均

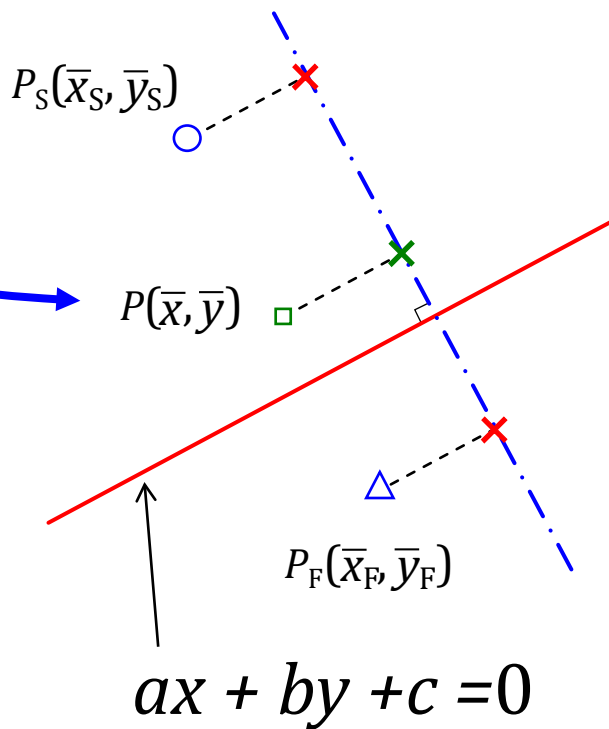
$$\bar{x} = \frac{1}{n} \left\{ \sum_{i=1}^{n_S} x_{Si} + \sum_{i=1}^{n_F} x_{Fi} \right\}$$

$$\bar{y} = \frac{1}{n} \left\{ \sum_{i=1}^{n_S} y_{Si} + \sum_{i=1}^{n_F} y_{Fi} \right\}$$

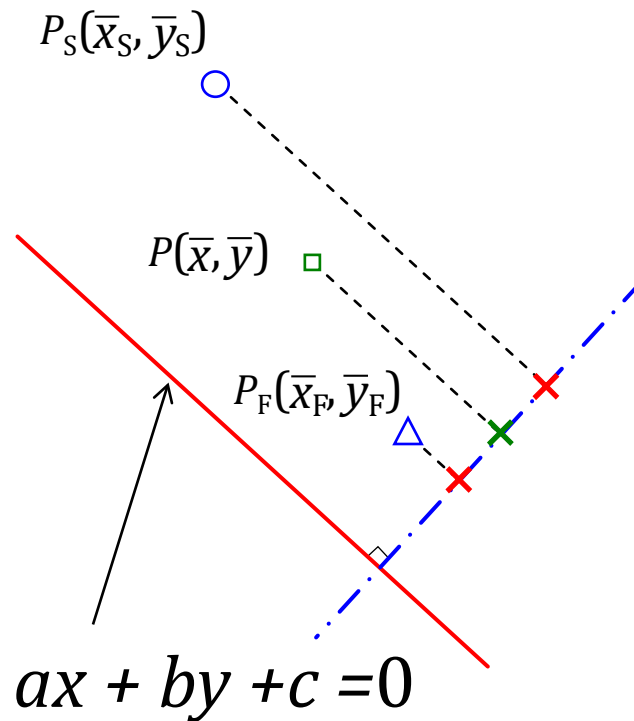
クラス間散布 s_B^2

$$= n_S(a\bar{x}_S + b\bar{y}_S - a\bar{x} - b\bar{y})^2$$

$$+ n_F(a\bar{x}_F + b\bar{y}_F - a\bar{x} - b\bar{y})^2$$

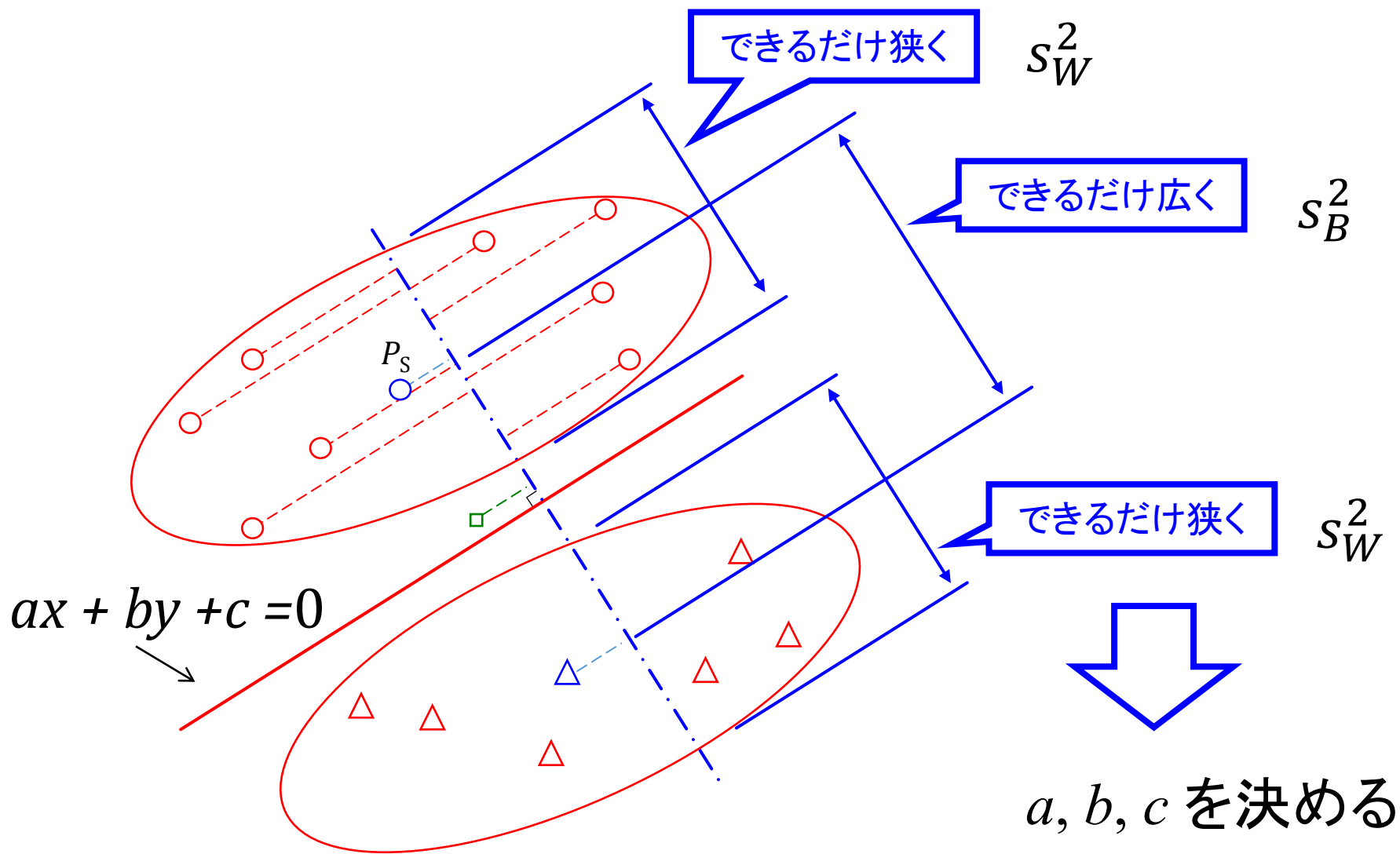


(a) クラス間散布大



(b) クラス間散布小

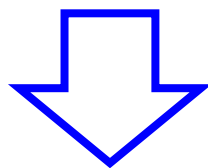
図2.6 判別線の傾きとクラス間散布



指標: $F = \frac{s_W^2}{s_B^2} \Rightarrow$ 最小

2.2 行列・ベクトルによる表現

$$F = \frac{S_W^2}{S_B^2} = \frac{\sum_{i=1}^{n_S} (ax_{Si} + by_{Si} - a\bar{x}_S - b\bar{y}_S)^2 + \sum_{i=1}^{n_F} (ax_{Fi} + by_{Fi} - a\bar{x}_F - b\bar{y}_F)^2}{n_S(a\bar{x}_S + b\bar{y}_S - a\bar{x} - b\bar{y})^2 + n_F(a\bar{x}_F + b\bar{y}_F - a\bar{x} - b\bar{y})^2}$$



行列, ベクトルにより表現する.

$$W = \begin{pmatrix} x_{S1} - \bar{x}_S & y_{S1} - \bar{y}_S \\ x_{S2} - \bar{x}_S & y_{S2} - \bar{y}_S \\ \vdots & \vdots \\ x_{Sn_S} - \bar{x}_S & y_{Sn_S} - \bar{y}_S \\ x_{F1} - \bar{x}_F & y_{F1} - \bar{y}_F \\ x_{F2} - \bar{x}_F & y_{F2} - \bar{y}_F \\ \vdots & \vdots \\ x_{Fn_F} - \bar{x}_F & y_{Fn_F} - \bar{y}_F \end{pmatrix}$$

$$\mathbf{v} = \begin{pmatrix} a \\ b \end{pmatrix}$$

$$\mathbf{Z} = W\mathbf{v} \quad (2.33)$$

$$\begin{aligned} S_W^2 &= \mathbf{Z}^t \mathbf{Z} \\ &= \mathbf{v}^t W^t W \mathbf{v} \\ &= \mathbf{v}^t U_W \mathbf{v} \end{aligned} \quad (2.34)$$

$$U_W = W^t W \quad (2.35)$$

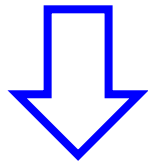
$$\mathbf{B} = \begin{pmatrix} \bar{x}_S - \bar{x} & \bar{y}_S - \bar{y} \\ \bar{x}_S - \bar{x} & \bar{y}_S - \bar{y} \\ \vdots & \vdots \\ \bar{x}_S - \bar{x} & \bar{y}_S - \bar{y} \\ \bar{x}_F - \bar{x} & \bar{y}_F - \bar{y} \\ \bar{x}_F - \bar{x} & \bar{y}_F - \bar{y} \\ \vdots & \vdots \\ \bar{x}_F - \bar{x} & \bar{y}_F - \bar{y} \end{pmatrix} \quad (2.36)$$

$$\begin{aligned}
 S_B^2 &= \mathbf{Y}^t \mathbf{Y} \\
 &= \mathbf{v}^t \mathbf{B}^t \mathbf{B} \mathbf{v} \\
 &= \mathbf{v}^t \mathbf{U}_B \mathbf{v}
 \end{aligned} \quad (2.37)$$

$$\mathbf{U}_B = \mathbf{B}^t \mathbf{B} \quad (2.38)$$

$$\mathbf{Y} = \mathbf{B} \mathbf{v}$$

$$F = \frac{S_W^2}{S_B^2} = \frac{\mathbf{v}^t \mathbf{U}_W \mathbf{v}}{\mathbf{v}^t \mathbf{U}_B \mathbf{v}} \quad (2.39)$$



F を最小とする \mathbf{v}

$$\mathbf{U}_W^{-1} \mathbf{U}_B \mathbf{v} = \lambda \mathbf{v} \quad (2.47)$$

$$\Rightarrow \mathbf{v} = \begin{pmatrix} a \\ b \end{pmatrix} \quad a, b \text{ の決定}$$

2.3 Rによる行列計算

判別分析_行列計算_入社試験の成績.R

```
x <- read.csv("C:/Users/Furuhashi/Documents/日本テクノセンター/判別分析  
/test_for_employment.csv")
```

表2.1 入社試験の成績と入社後の評価

学科	面接	10年後の評価
6	6	1
3	6	1
3	3	1
6	3	2
6	0	2
3	0	2

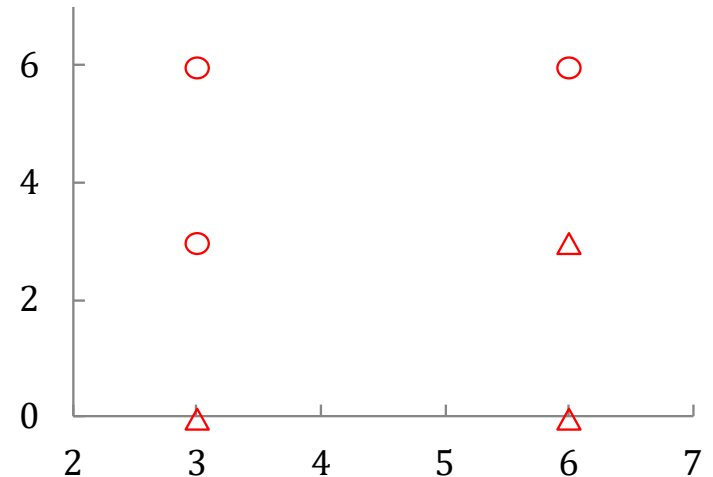


図2.7 入社試験の成績と入社後の評価

```
x_S <- subset(x[,-3], x$入社後の評価 < 1.5)
```

```
x_F <- subset(x[,-3], x$入社後の評価 > 1.5)
```

学科 面接

1	6	6
2	3	6
3	3	3

$$\mathbf{x}_S = \begin{pmatrix} x_{S1} & y_{S1} \\ x_{S2} & y_{S2} \\ x_{S3} & y_{S3} \end{pmatrix}$$

学科 面接

4	6	3
5	6	0
6	3	0

$$\mathbf{x}_F = \begin{pmatrix} x_{S1} & y_{S1} \\ x_{S2} & y_{S2} \\ x_{S3} & y_{S3} \end{pmatrix}$$

```
WS <- scale(x_S,center = TRUE, scale = FALSE)
```

```
WF <- scale(x_F,center = TRUE, scale = FALSE)
```

学科 面接

1	2	1
2	-1	1
3	-1	-2

$$\mathbf{W}_S = \begin{pmatrix} x_{S1} - \bar{x}_S & y_{S1} - \bar{y}_S \\ x_{S2} - \bar{x}_S & y_{S2} - \bar{y}_S \\ x_{S3} - \bar{x}_S & y_{S3} - \bar{y}_S \end{pmatrix}$$

学科 面接

4	1	2
5	1	-1
6	-2	-1

$$\mathbf{W}_F = \begin{pmatrix} x_{S1} - \bar{x}_S & y_{S1} - \bar{y}_S \\ x_{S2} - \bar{x}_S & y_{S2} - \bar{y}_S \\ x_{S3} - \bar{x}_S & y_{S3} - \bar{y}_S \end{pmatrix}$$

```
UW <- t(WS) %*% WS + t(WF) %*% WF
```

$$\mathbf{U}_W = \mathbf{W}^t \mathbf{W}$$

```
UB <- num_S*((mean_S-mean_全体) %*% t(mean_S-mean_全体)) +  
      num_F*((mean_F-mean_全体) %*% t(mean_F-mean_全体))
```

$$\mathbf{U}_B = \mathbf{B}^t \mathbf{B}$$

```
inv_UW <- solve(UW)
```

$$U_W^{-1}$$

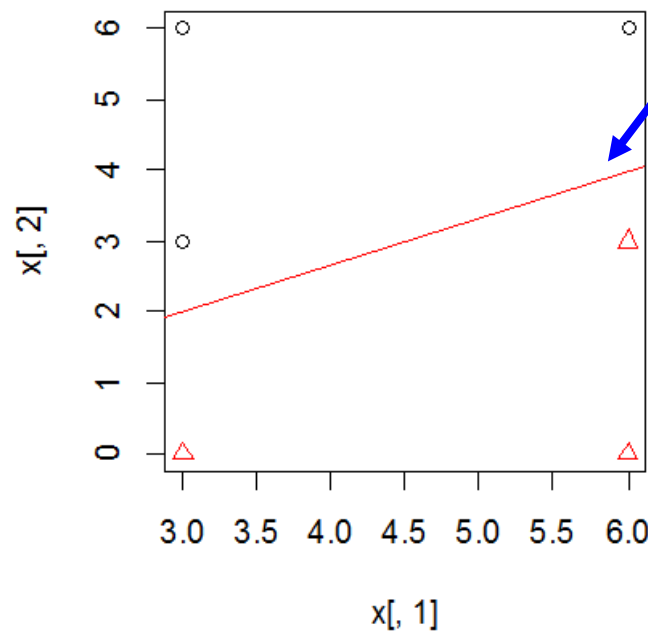
```
Eigen_value <- eigen(inv_UW %*% UB) a, b の決定
```

➡ $U_W^{-1}U_B$ の固有値, 固有ベクトル計算 $v = \begin{pmatrix} a \\ b \end{pmatrix}$

```
abline(0,  $-\frac{\text{Re}(\text{Eigen\_value}\$vectors[1,1])}{\text{Re}(\text{Eigen\_value}\$vectors[2,1])}$ , col = "red")
```

$-c/b$

$-a/b$



$$ax + by + c = 0$$

$$y = -\frac{a}{b}x - \frac{c}{b}$$

図2.11 abline()の実行結果

3. 定数項cの決め方

3.1 マハラノビス距離

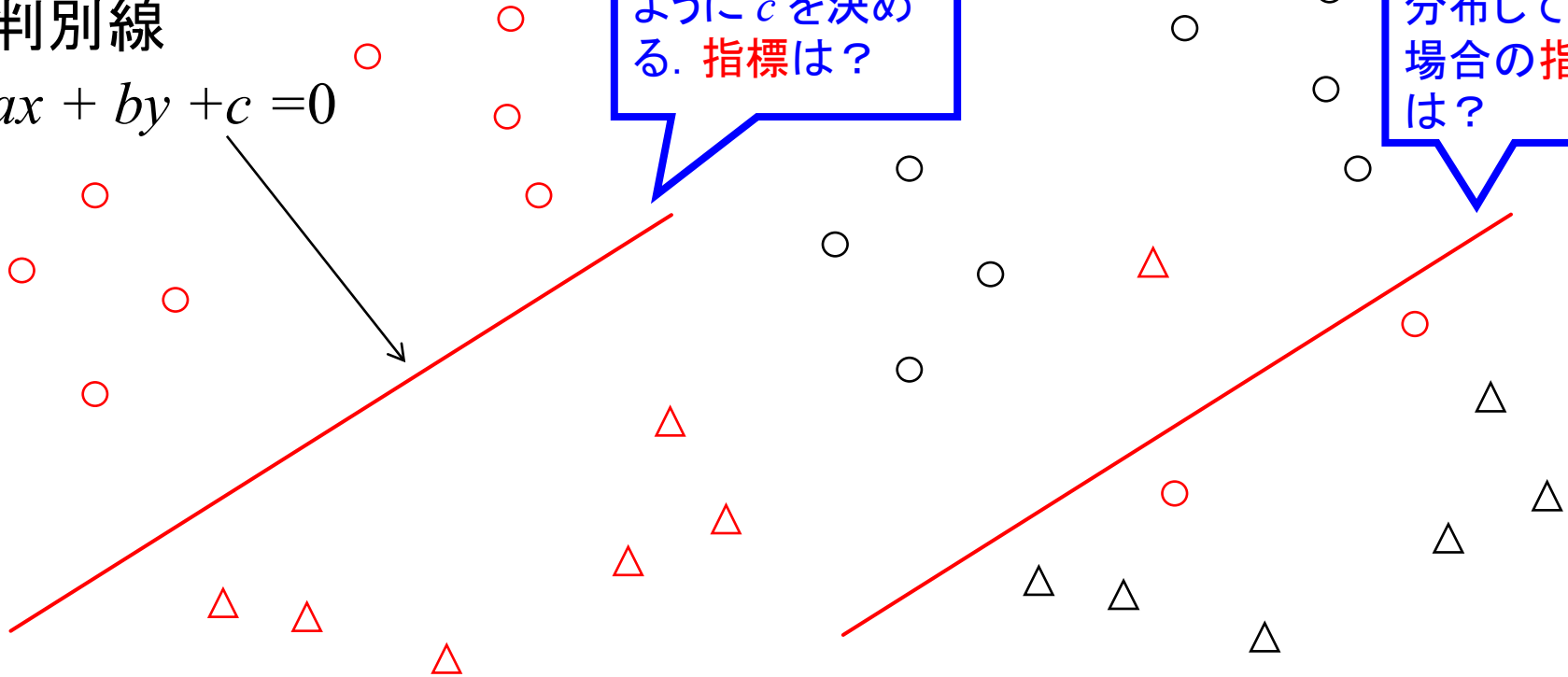
判別線の定数項cはどう決めればよいか？

判別線

$$ax + by + c = 0$$

両グループが
できるだけ離れる
ようにcを決め
る. 指標は？

両グループが
混じり合って
分布している
場合の指標
は？



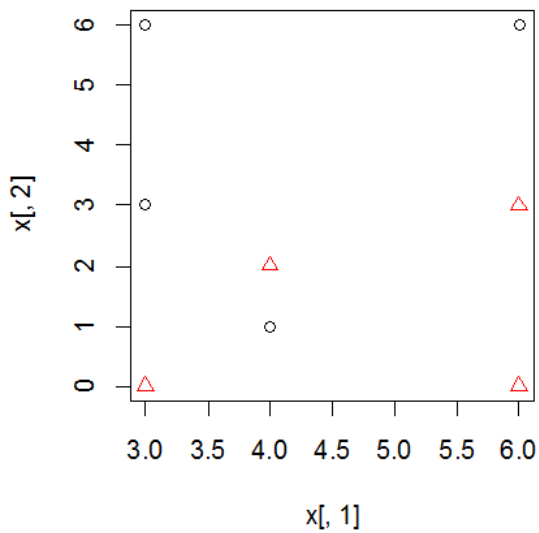


図3.8 表3.1のデータ分布の例

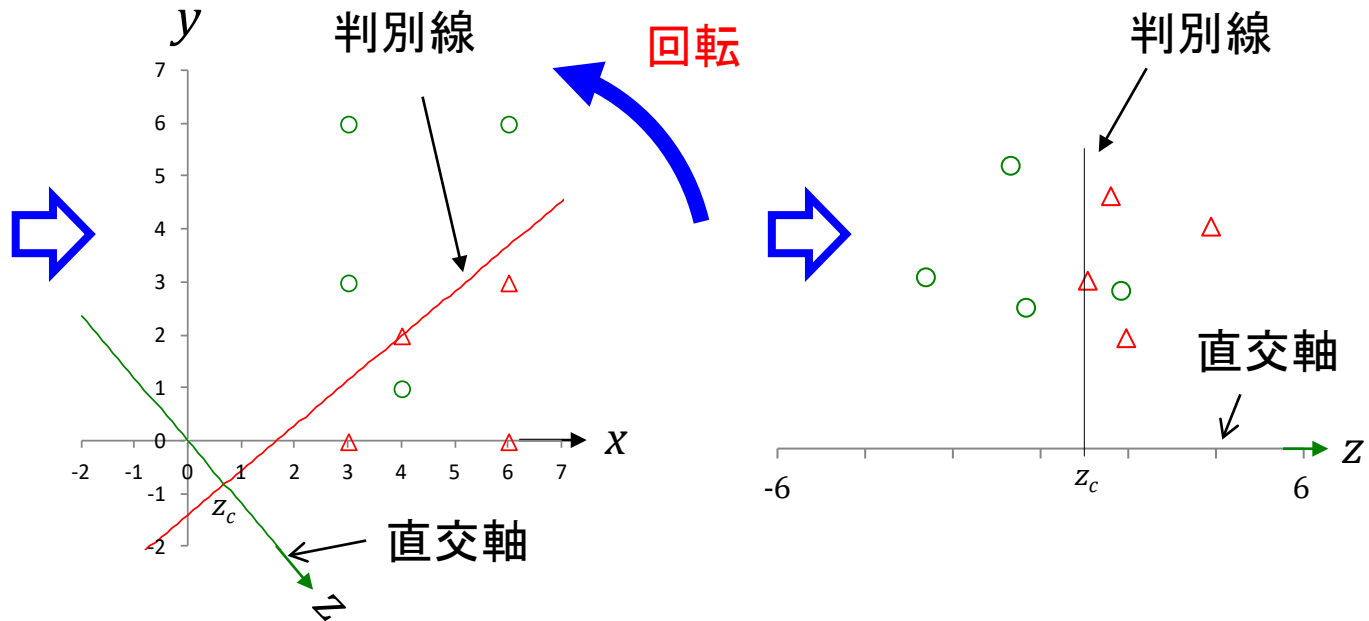
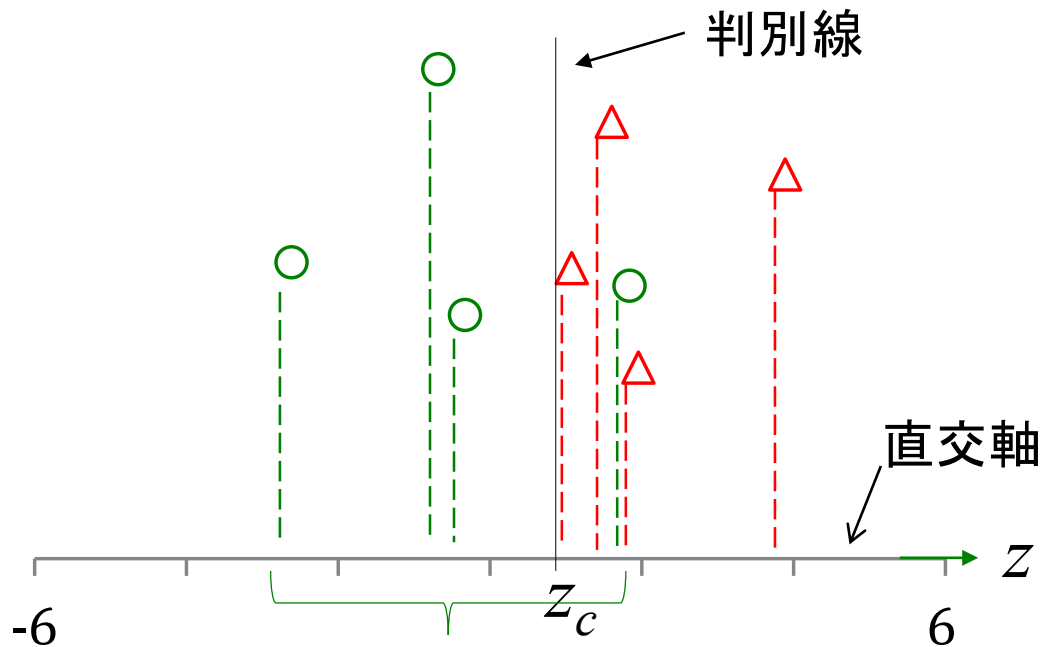


図3.10 判別線と直交軸 z

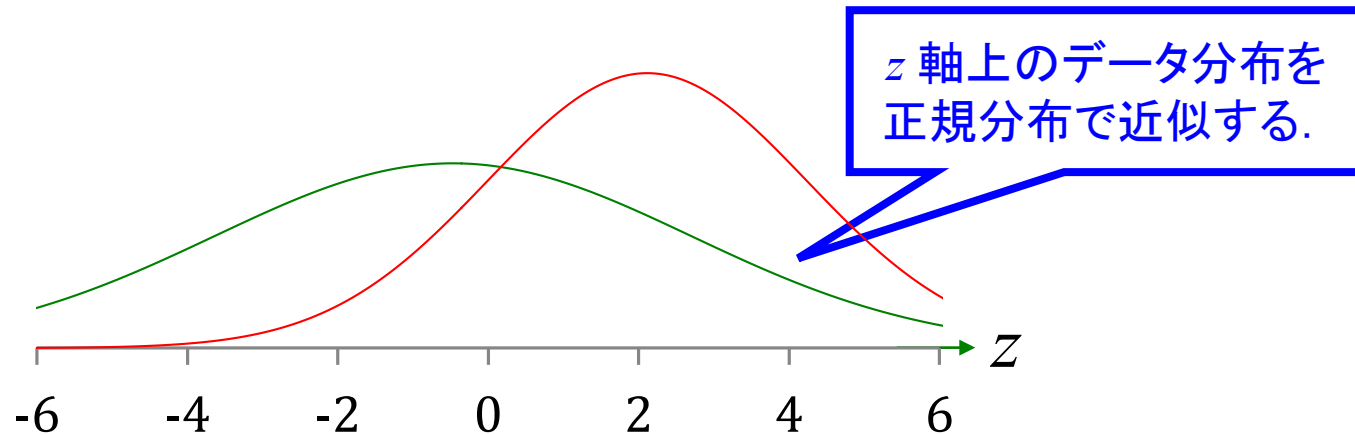


平均 $\bar{z}_S = a\bar{x}_S + b\bar{y}_S + c_0$

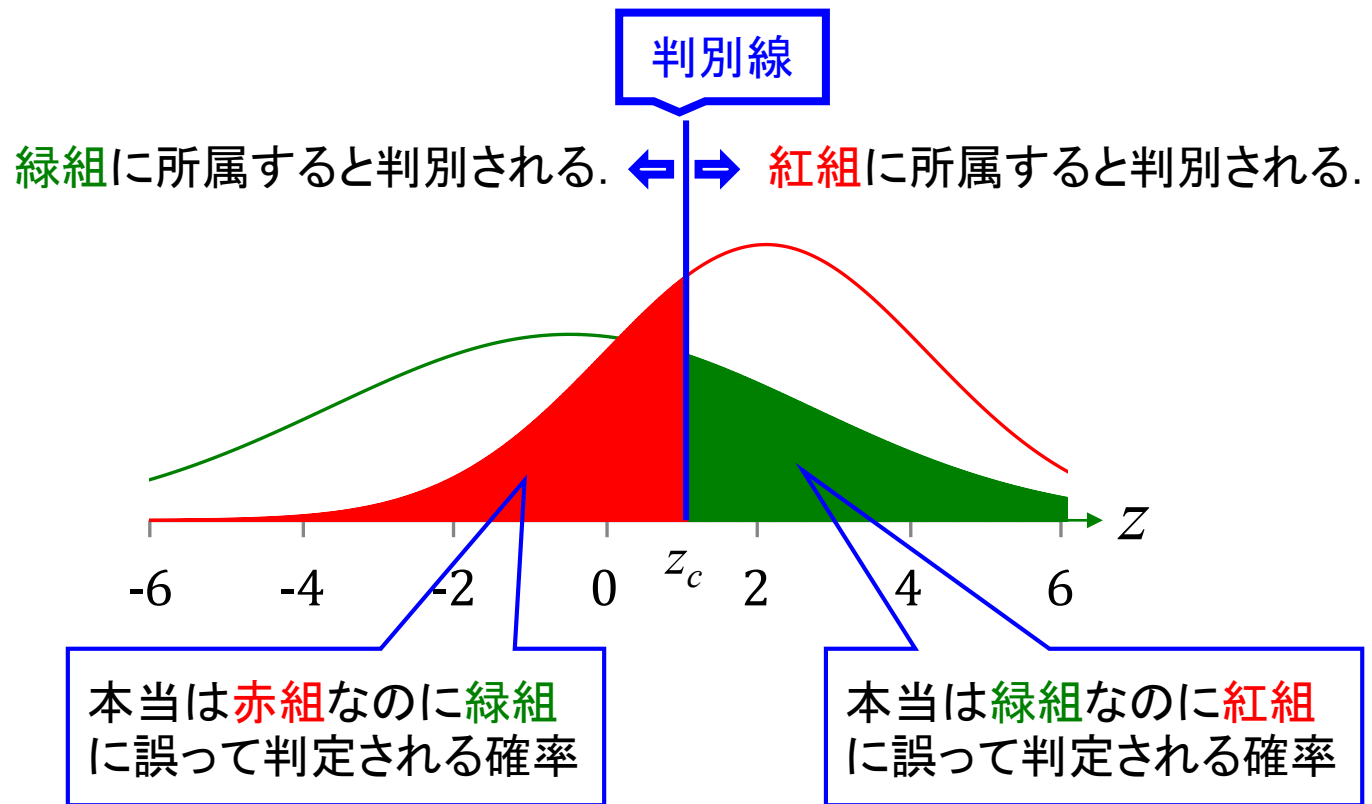
不偏分散 $v_{eS}^2 = \frac{1}{n_S - 1} \sum_{i=1}^{n_S} (z_{Si} - \bar{z}_S)^2$



$$f_S(z) = \frac{1}{\sqrt{2\pi v_{eS}^2}} e^{\left\{ -\frac{(z - \bar{z}_S)^2}{2v_{eS}^2} \right\}}$$



新しいデータが来たときにそのデータの所属先を決めるには？



z_c はそれぞれの間違い確率を同じにする.

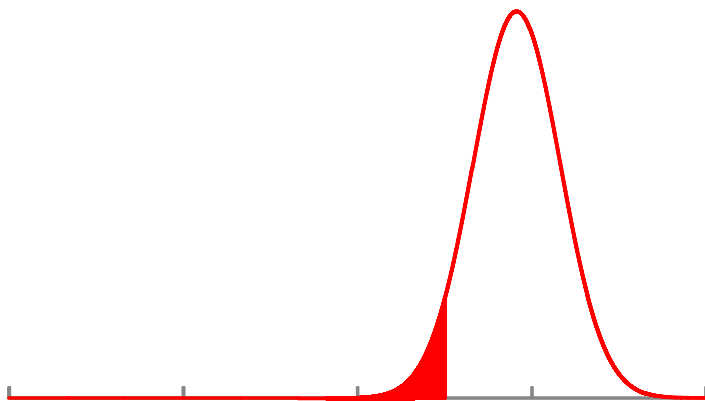
緑と赤の塗りつぶした箇所の面積を一緒にする z_c は

マハラノビス距離:
$$\frac{z_c - \bar{z}_S}{\sqrt{v_{eS}^2}} = \frac{-z_c + \bar{z}_F}{\sqrt{v_{eF}^2}}$$

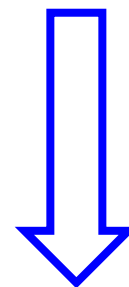
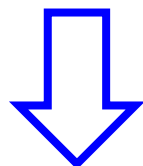
 z_c の決定

$$z_c = \frac{\sqrt{v_{eS}^2} \bar{z}_F + \sqrt{v_{eF}^2} \bar{z}_S}{\sqrt{v_{eS}^2} + \sqrt{v_{eF}^2}}$$

マハラノビス距離より大きな値をとる確率が同じである理由



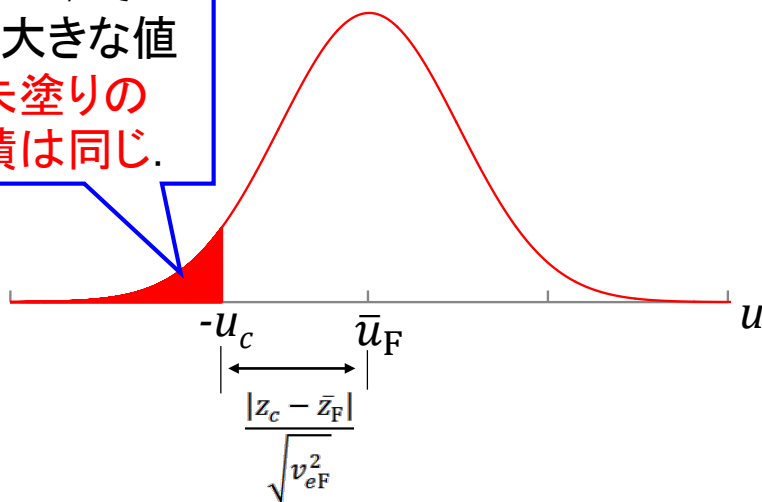
$$f_F(z) = \frac{1}{\sqrt{2\pi v_{eF}^2}} e^{-\frac{(z - \bar{z}_F)^2}{2v_{eF}^2}}$$



マハラノビス距離が同じであれば、それより大きな値の朱塗りの面積は同じ.

$$u = \frac{z - \bar{z}_F}{v_{eF}} \quad \text{マハラノビス距離の変数変換により}$$

$$dz = v_{eF} du$$



平均0, 分散1の標準分布に変換される.

$$f_F(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}}$$

3.3 Rによる計算

判別分析_定数項の決定_入社試験の成績.R

表3.1 入社試験の成績と入社後の評価(II)

学科	面接	入社後の評価
6	6	1
3	6	1
3	3	1
6	3	2
6	0	2
3	0	2
4	1	1
4	2	2

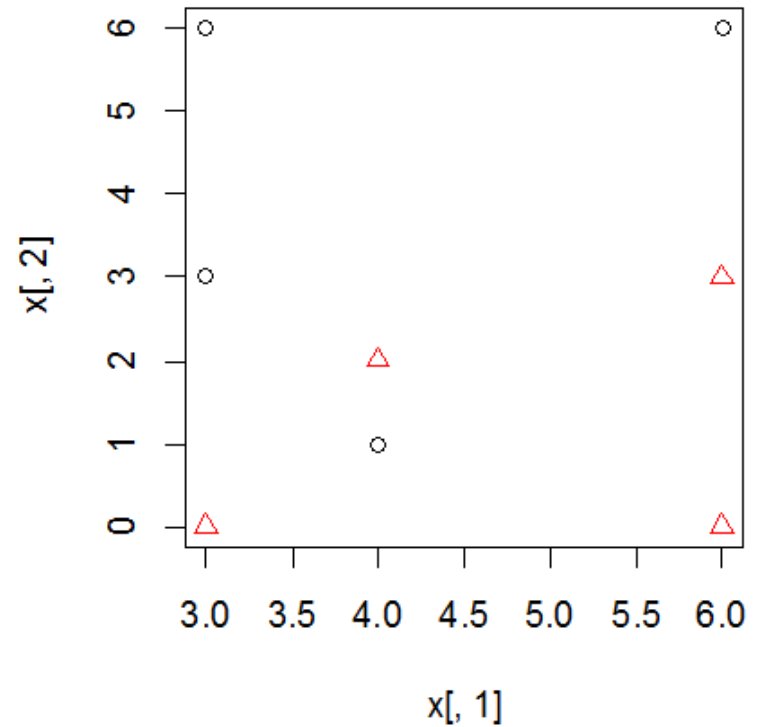


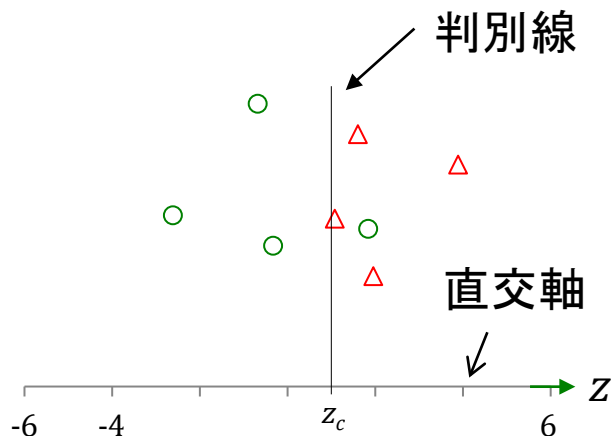
図3.8 表3.1のデータ分布

```
ve_S <- (t(v) %*% t(WS) %*% WS %*% v)/(num_S-1) #veS^2の計算
```

```
ve_F <- (t(v) %*% t(WF) %*% WF %*% v)/(num_F-1) #veF^2の計算
```

```
mean_zS <- t(v) %*% mean_S #zSの平均値の計算(c = 0の仮定の下で)
```

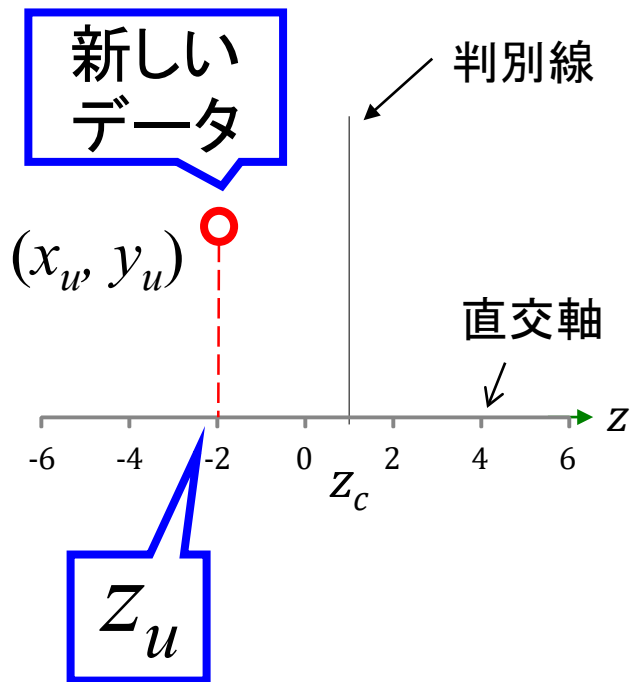
```
mean_zF <- t(v) %*% mean_F #zFの平均値の計算(c = 0の仮定の下で)
```



```
zc <- (sqrt(ve_F)*mean_zS + sqrt(ve_S)*mean_zF)/(sqrt(ve_F) + sqrt(ve_S))
```

#zcの計算

3.3 未知データの判別



$$z_u = ax_u + by_u$$

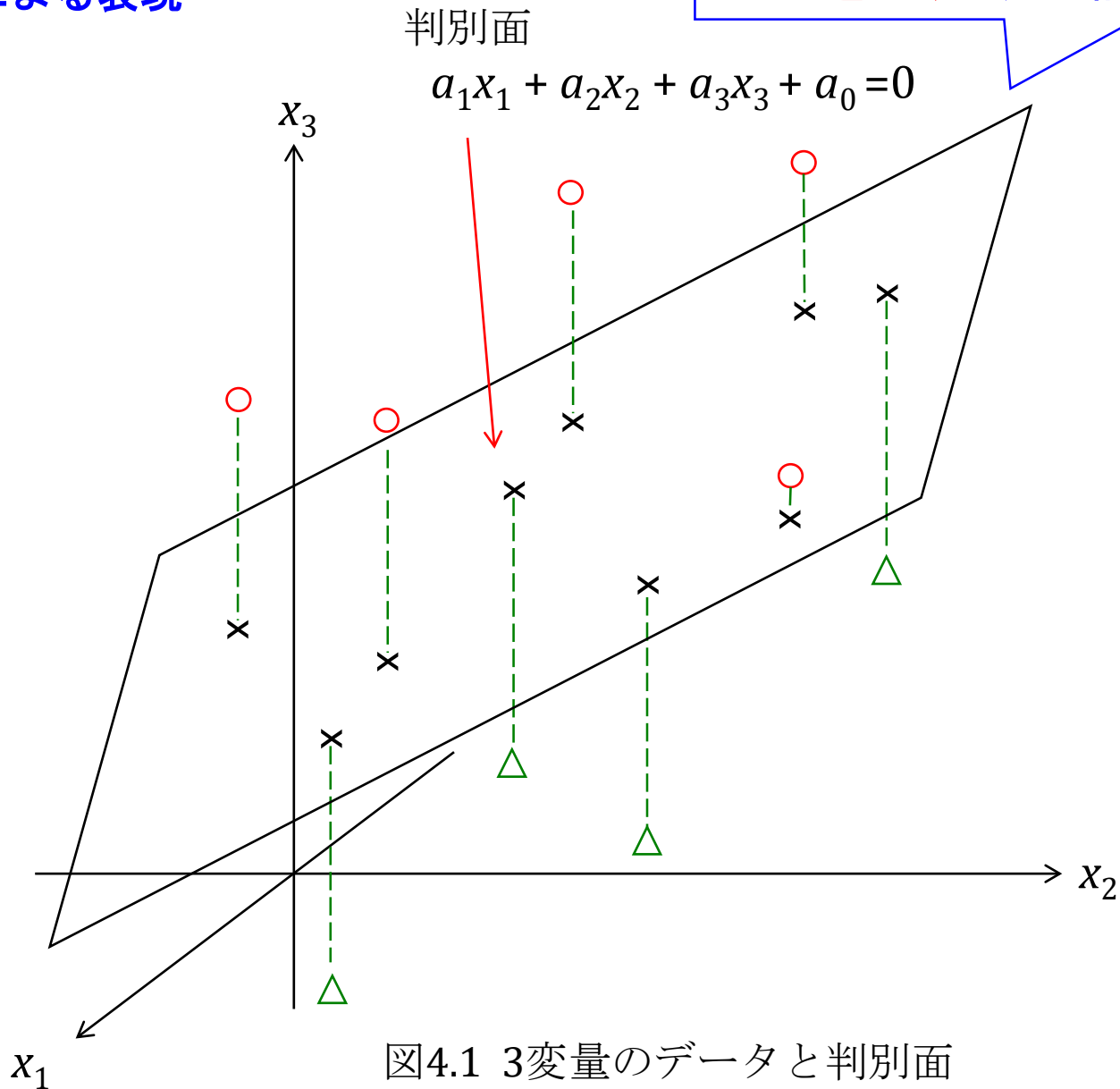
$z_u > z_c$: 低評価

$z_u < z_c$: 高評価

4. 線形判別分析 (多変量)

4.1 行列による表現

二つのクラスを分ける面を同定する課題



$$W = \begin{pmatrix} x_{S11} - \bar{x}_{S1} & x_{S21} - \bar{x}_{S2} & \cdots & x_{Sk1} - \bar{x}_{Sk} \\ x_{S12} - \bar{x}_{S1} & x_{S22} - \bar{x}_{S2} & \cdots & x_{Sk2} - \bar{x}_{Sk} \\ \vdots & \vdots & \ddots & \vdots \\ x_{S1n_S} - \bar{x}_{S1} & x_{S2n_S} - \bar{x}_{S2} & \cdots & x_{Skn_S} - \bar{x}_{Sk} \\ x_{F11} - \bar{x}_{F1} & x_{F21} - \bar{x}_{F2} & \cdots & x_{Fk1} - \bar{x}_{Fk} \\ x_{F12} - \bar{x}_{F1} & x_{F22} - \bar{x}_{F2} & \cdots & x_{Fk2} - \bar{x}_{Fk} \\ \vdots & \vdots & \ddots & \vdots \\ x_{F1n_F} - \bar{x}_{F1} & x_{F2n_F} - \bar{x}_{F2} & \cdots & x_{Fkn_F} - \bar{x}_{Fk} \end{pmatrix} \quad (4.2)$$

$$v = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_k \end{pmatrix} \quad (4.3)$$

$$U_W = W^t W \quad (4.4)$$

$$B = \begin{pmatrix} \bar{x}_{S1} - \bar{x} & \bar{x}_{S2} - \bar{x} & \cdots & \bar{x}_{Sk} - \bar{x} \\ \bar{x}_{S1} - \bar{x} & \bar{x}_{S2} - \bar{x} & \cdots & \bar{x}_{Sk} - \bar{x} \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_{S1} - \bar{x} & \bar{x}_{S2} - \bar{x} & \cdots & \bar{x}_{Sk} - \bar{x} \\ \bar{x}_{F1} - \bar{x} & \bar{x}_{F2} - \bar{x} & \cdots & \bar{x}_{Fk} - \bar{x} \\ \bar{x}_{F1} - \bar{x} & \bar{x}_{F2} - \bar{x} & \cdots & \bar{x}_{Fk} - \bar{x} \\ \vdots & \vdots & \vdots & \vdots \\ \bar{x}_{F1} - \bar{x} & \bar{x}_{F2} - \bar{x} & \cdots & \bar{x}_{Fk} - \bar{x} \end{pmatrix} \quad (4.5)$$

$$U_B = B^t B \quad (4.6)$$

$$U_W^{-1} U_B v = \lambda v \quad (4.7)$$

次数が増えても表式は同じ

4.2 定数項 a_0 の決定

仮の判別面($a_0 = 0$)

$$a_1x_1 + a_2x_2 + a_3x_3 = 0$$

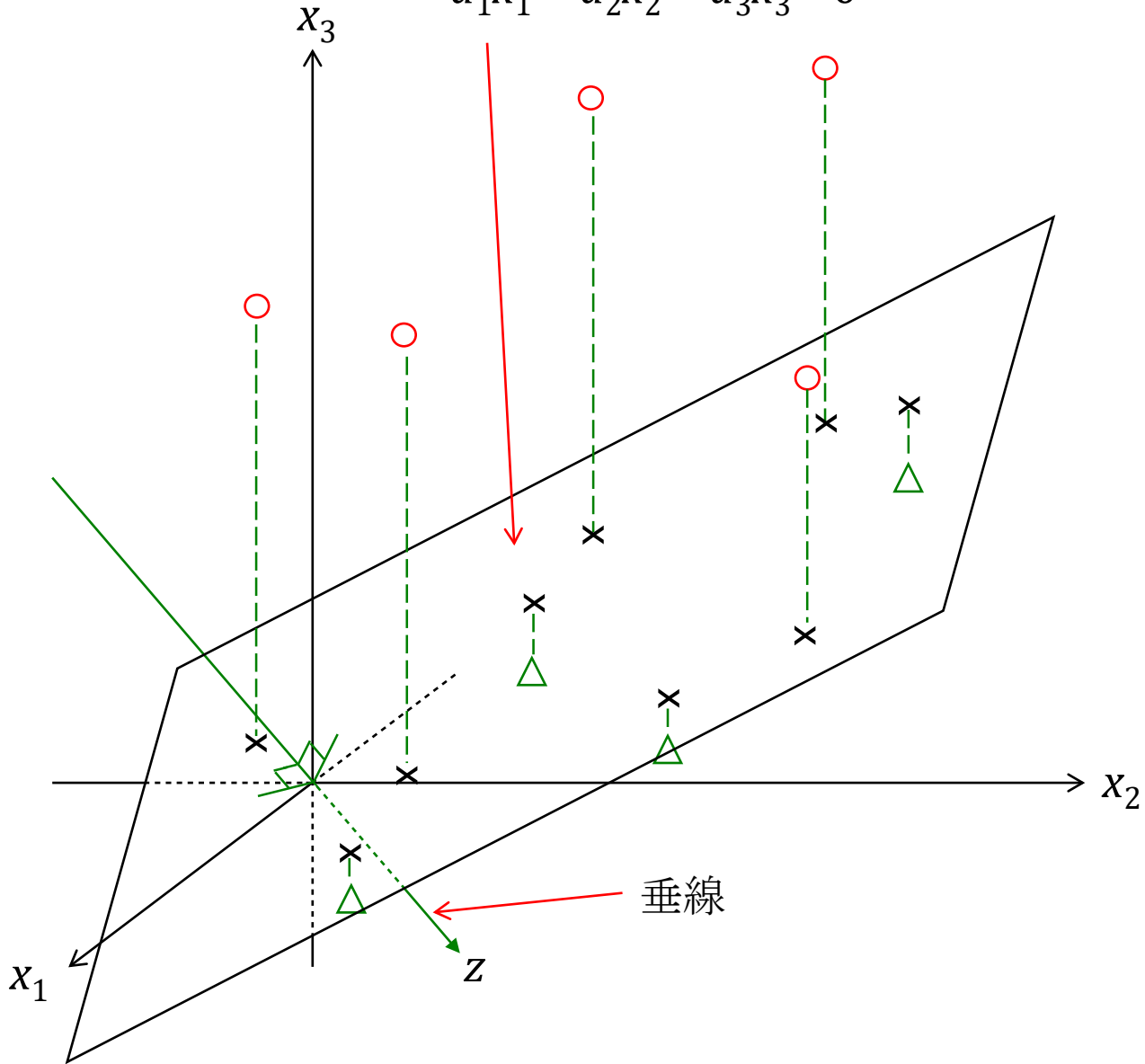


図4.2 仮の判別面($a_0 = 0$)と垂線

判別面 : $a_1x_1 + a_2x_2 + a_3x_3 + a_0 = 0$

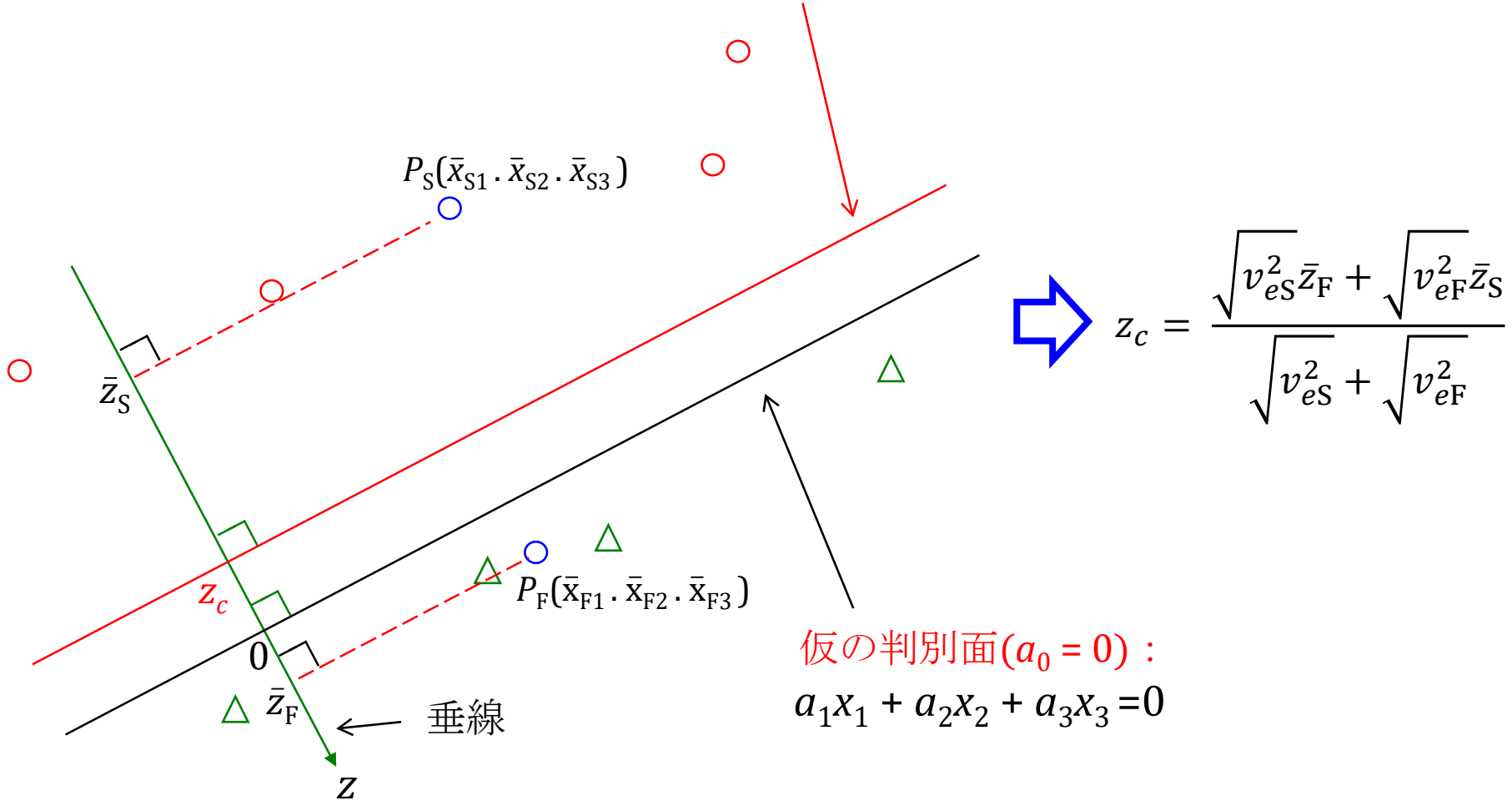


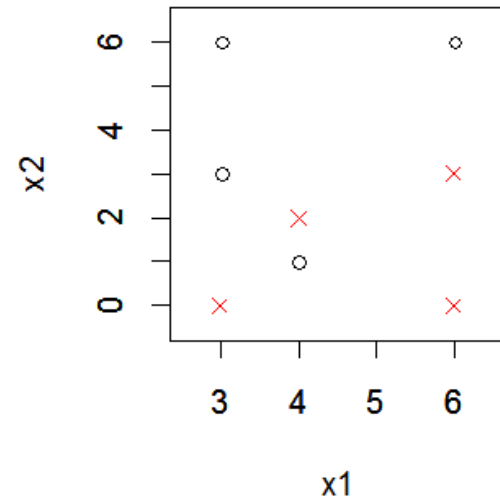
図4.3 仮の判別面の垂線上への射影

4.3 Rによる行列計算

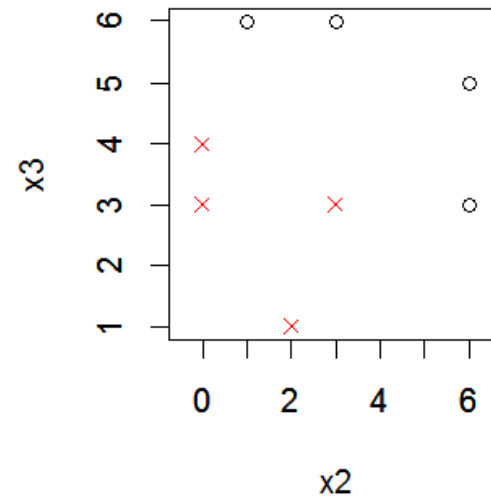
判別分析_3変量_入社試験の成績.R

表4.1 入社試験の成績と入社後の評価
(3変量)

学科	面接	実技	入社後の評価
6	6	5	1
3	6	3	1
3	3	6	1
6	3	3	2
6	0	4	2
3	0	3	2
4	1	6	1
4	2	1	2



(a) $x_1 - x_2$



(b) $x_2 - x_3$

図4.4 表4.1のデータ分布

```

UB <- num_S*((mean_S-mean_全体) %*% t(mean_S-mean_全体)) +
      num_F*((mean_F-mean_全体) %*% t(mean_F-mean_全体))
      # UB行列を求める
UW <- t(WS) %*% WS + t(WF) %*% WF
      # UW行列を求める
inv_UW <- solve(UW)
      # UW行列の逆行列を求める
Eigen_value <- eigen(inv_UW %*% UB)
      # UW^(-1) UB の固有値計算
v <- Re(Eigen_value$vector[,1])
      #固有ベクトルをvに代入

```

判別分析

```

ve_S <- (t(v) %*% t(WS) %*% WS %*% v)/(num_S-1)
      #veS^2の計算
ve_F <- (t(v) %*% t(WF) %*% WF %*% v)/(num_F-1)
      #veF^2の計算
mean_zS <- t(v) %*% mean_S
      #zSの平均値の計算(c = 0)
mean_zF <- t(v) %*% mean_F
      #zFの平均値の計算(c = 0)

zc <- (sqrt(ve_F)*mean_zS + sqrt(ve_S)*mean_zF)/(sqrt(ve_F) + sqrt(ve_S))
      #zcの計算

```

定数項 の計算

```

x1_cd <- seq(min_x1, max_x1, length = 50)
      # x1軸の[min, max]区間を50分割座標の生成
x2_cd <- seq(min_x2, max_x2, length = 50)
      # x2軸の[min, max]区間を50分割座標の生成
x1 <- matrix(0, 2500)
x2 <- matrix(0, 2500)
y <- matrix(0, 2500)
for(i in 1:50) for(j in 1:50){x1[(i-1)*50+j] <- x1_cd[i];
x2[(i-1)*50+j] <- x2_cd[j]; y[(i-1)*50 + j] <- x3_intercept - v[1]/v[3]*x1_cd[i] -
v[2]/v[3]*x2_cd[j]}
      # x1-x2平面の50*50 = 2500の格子点上のyの値の計算
plot3d(x1, x2, y, col="blue", pch=2, add = 1)

```

判別面 の描画

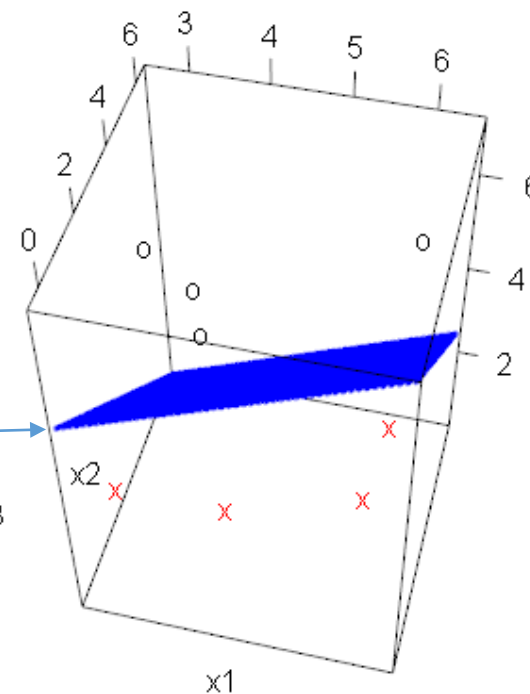


図4.5 判別分析結果

5. 非線形判別分析

5.2 マハラノビス距離

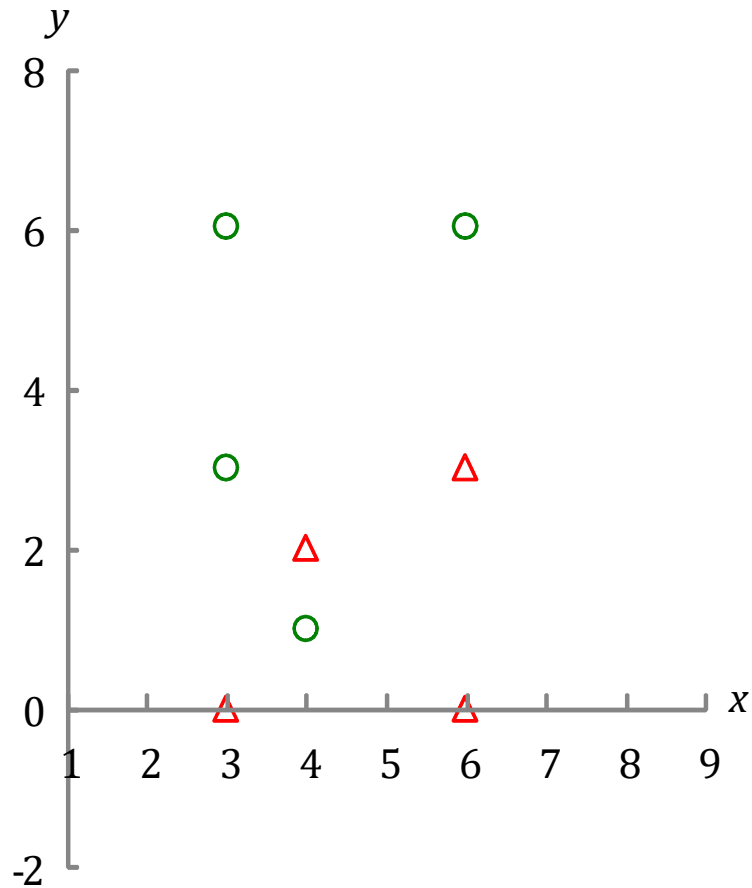


図5.12 入社試験の成績

判別面をどのように引けばそれぞれに間違う確率を同じにできるか？

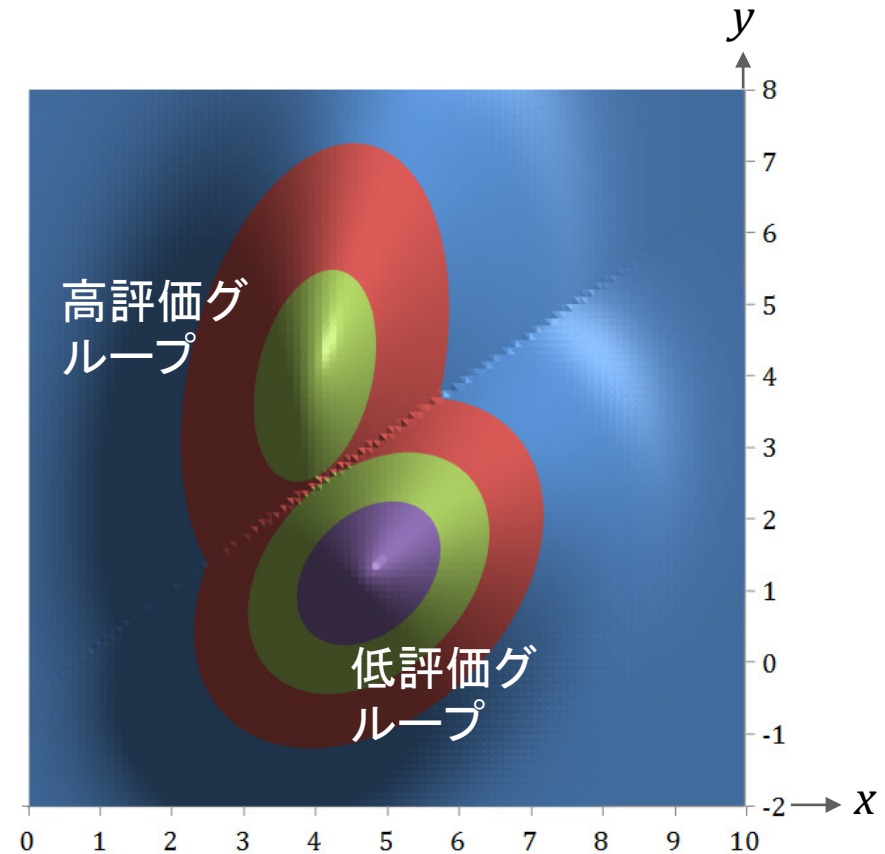


図5.9 2変量正規分布でデータ分布を近似

$$x_S = 4, \quad y_S = 4, \quad v_{xS}^2 = 2, \quad v_{yS}^2 = 6, \quad \rho = 0.289$$
$$x_F = 4.75, \quad y_F = 1.25, \quad v_{xF}^2 = 2.25, \quad v_{yF}^2 = 1.25, \quad \rho = 0.333$$

2変量を x, y とし, クラス S に属するデータの平均 \bar{x}_S, \bar{y}_S と不偏分散 v_{xS}^2, v_{yS}^2 を

$$\bar{x}_S = \frac{1}{n_S} \sum_{i=1}^{n_S} x_{Si}, \quad \bar{y}_S = \frac{1}{n_S} \sum_{i=1}^{n_S} y_{Si} \quad (5.1)$$

$$v_{xS}^2 = \frac{1}{n_S - 1} \sum_{i=1}^{n_S} (x_{Si} - \bar{x}_S)^2 \quad (5.2)$$

$$v_{yS}^2 = \frac{1}{n_S - 1} \sum_{i=1}^{n_S} (y_{Si} - \bar{y}_S)^2$$

とすると

$$f_S(x, y) = \frac{1}{2\pi v_{xS} v_{yS} \sqrt{(1 - \rho^2)}} \times \exp \left\{ -\frac{1}{2(1 - \rho^2)} \left\{ \frac{(x - \bar{x}_S)^2}{v_{xS}^2} - 2\rho \frac{(x - \bar{x}_S)(y - \bar{y}_S)}{v_{xS} v_{yS}} + \frac{(y - \bar{y}_S)^2}{v_{yS}^2} \right\} \right\} \quad (5.3)$$

ρ は変量 x, y 間の相関係数

$$\rho = \frac{v_{xyS}}{v_{xS} v_{yS}} \quad (5.4) \quad \text{ただし, } v_{xyS} = \frac{1}{n_S - 1} \sum_{i=1}^{n_S} (x_{Si} - \bar{x}_S)(y_{Si} - \bar{y}_S)$$

5.4 2変量正規分布の行列表現

分散共分散行列 Σ を

$$\Sigma_S = \begin{pmatrix} v_{xS}^2 & v_{xyS} \\ v_{xyS} & v_{yS}^2 \end{pmatrix} \quad (5.18)$$
$$\mathbf{X}_S = \begin{pmatrix} x - \bar{x}_S \\ y - \bar{y}_S \end{pmatrix}$$

と定義すると, (5.3)式の2次元正規分布 $f_S(x, y)$ は

$$f_S(x, y) = \frac{1}{2\pi\sqrt{|\Sigma_S|}} \exp\left\{-\frac{1}{2} \mathbf{X}_S^t \Sigma_S^{-1} \mathbf{X}_S\right\} \quad (5.19)$$

マハラノビス距離
の2乗値

ただし, $|\Sigma|$ は行列 Σ の行列式, Σ^{-1} は行列 Σ の逆行列

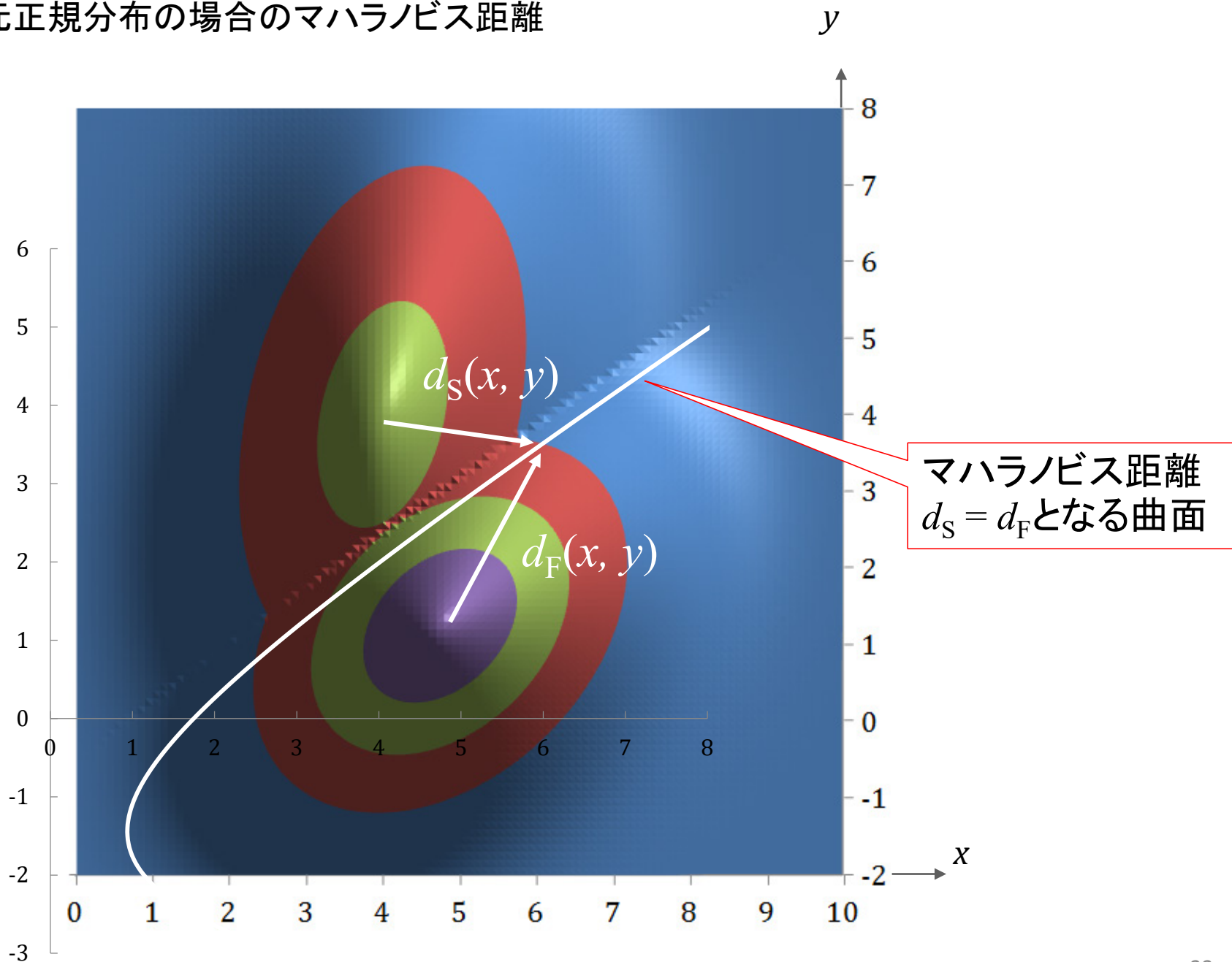
$$|\Sigma_S| = v_{xS}^2 v_{yS}^2 - v_{xyS}^2 \quad (5.20)$$

$$\Sigma_S^{-1} = \frac{1}{v_{xS}^2 v_{yS}^2 - v_{xyS}^2} \begin{pmatrix} v_{yS}^2 & -v_{xyS} \\ -v_{xyS} & v_{xS}^2 \end{pmatrix} \quad (5.21)$$

マハラノビス距離 $d_S(x, y)$

$$d_S(x, y)^2 = \mathbf{X}_S^t \Sigma_S^{-1} \mathbf{X}_S \quad (5.22)$$

2次元正規分布の場合のマハラノビス距離



判別分析_マハラノビス距離による判別_行列表現_入社試験の成績.R

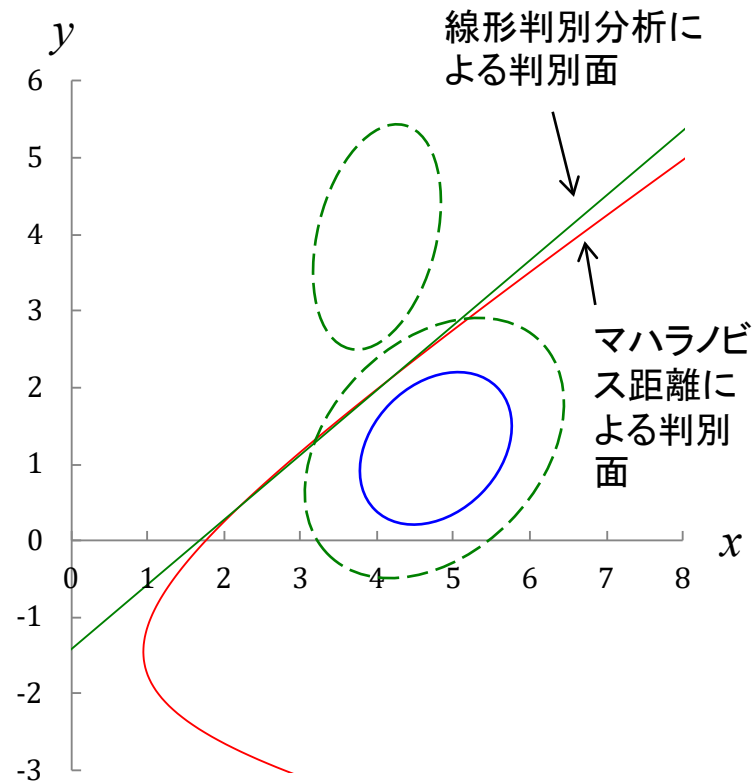
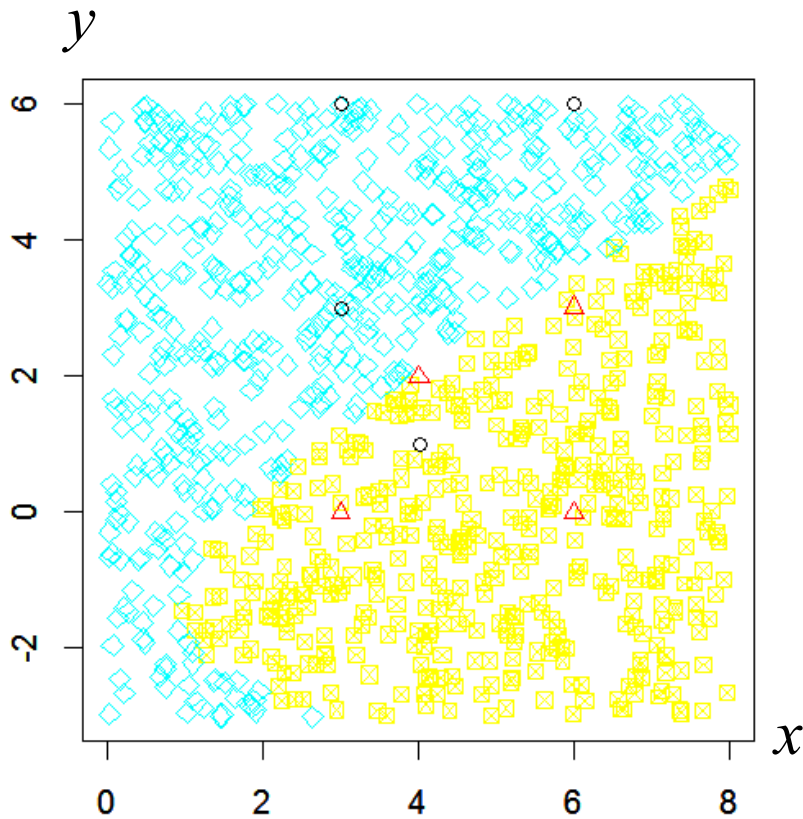


図5.15 入社試験の成績とマハラノビス距離による**判別結果**

図5.14 マハラノビス距離による判別面と線形判別分析による判別面の比較

5.7 Rによるアヤメのデータの判別

判別分析_マハラノビス距離による判別_行列表現_アイリス.R

```
xx <- data.frame(iris[,-5], as.integer(iris[,5]))
```

Rにはirisという名前でアヤメのデータが組み込まれている。

データの5列目はアヤメの種類

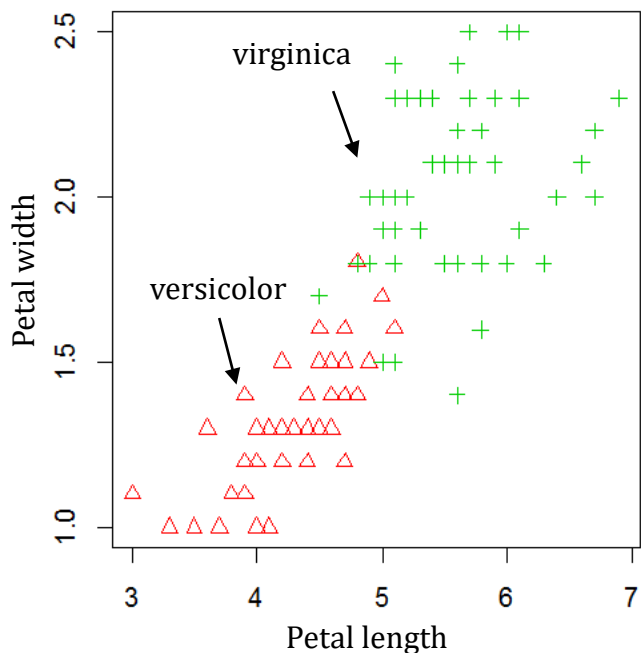


図5.16 アヤメのデータ
(versicolorとvirginica)

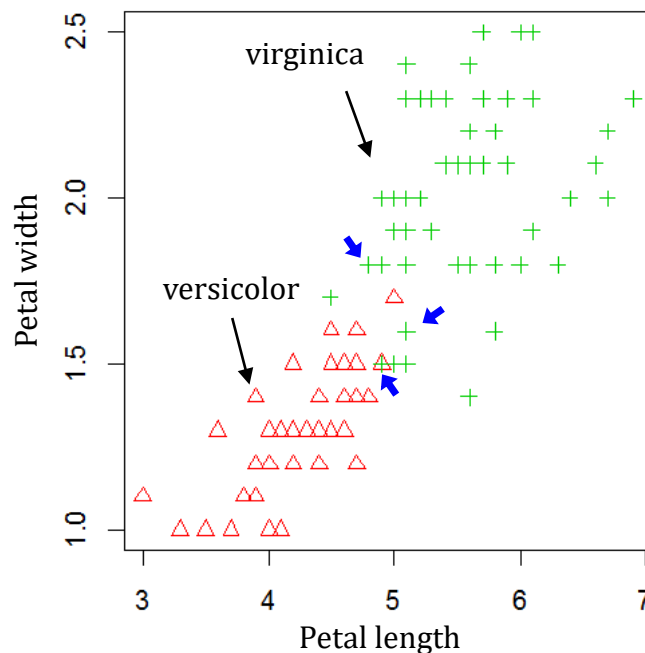


図5.17 アヤメのデータのマハラノ
ビス距離による再判別結果

判別分析_Ida_アイリス.R

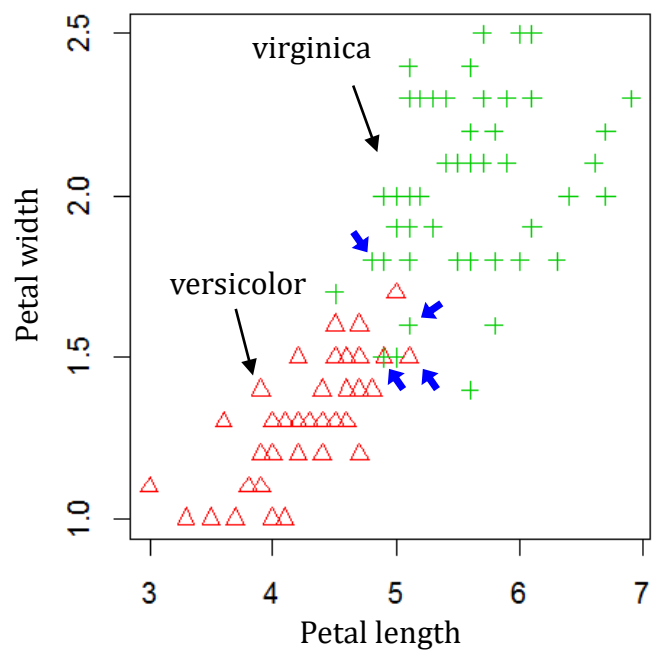


図5.18 アヤメのデータの線形判別
分析とマハラノビス距離に基づ
く z_c 決定法による再判別結果

本稿の詳細な解説は

[多変量解析の基礎 Ⅲ \(判別分析\) ー理論とRによる演習ー](#)

と題してアマゾンよりkindle版として出版しています。