

3. 判別分析 [【動画】](#)

重回帰分析では量的データの目的変数を複数の説明変数で直接予測するものですが、判別分析は質的な目的変数の分類を複数の説明変数で予測する手法です。質的データにはその値自身に意味はありませんので、直接値を予測するのではなく、判別関数という1次式を求め、その値によって分類を予測します。その際、予測の精度やどの変数の影響が強いかなどが検討されます。

今、試験の可否をそれまでの平均的な勉強時間と模擬試験の成績で判定するとします。合格、不合格という2群の場合の判別関数は、以下のように与えられます。

判別関数  $z = b_1 \text{ 勉強時間} + b_2 \text{ 平均点} + b_0$   
ここに、 $b_1, b_2, b_0$  は分析結果として決まるパラメータです。この判別関数の値が0以上か0未満かで可否を判別します。以下の図を見て下さい。

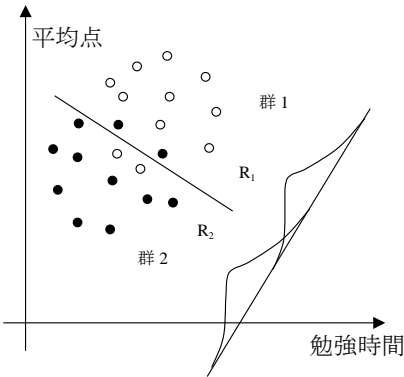


図1 判別分析での群分け

白丸を合格群、黒丸を不合格群として、その間に線を引き、その線より上を合格、下を不合格と判定することを考えます。その線が判別関数  $z = b_1 \text{ 勉強時間} + b_2 \text{ 平均点} + b_0 = 0$  を表しています。これを理論的にうまく引いてやるのが大切です。実際の分析を、例を使ってパソコンで見てみましょう。

例

入学試験の可否と勉強時間・模擬試験の平均点のデータを求めたところ以下のような結果を得た(判別分析.txt(p1))。可否を判定するための勉強時間と平均点の1次関数を求めよ。またこの関数によってデータを判別し、誤判別の確率を求めよ。

可否	勉強時間	平均点	可否	勉強時間	平均点
1	5.6	70.2	2	3.8	67.4
1	5.9	74.2	2	3.8	61.3
1	4.1	72.7	2	1.7	60.6
1	5.1	84.9	2	2.7	77.2
1	5.0	93.0	2	4.3	65.9
:	:	:	:	:	:
1	3.6	85.5	2	2.5	64.4
2	3.8	47.9	2	5.2	50.7
2	3.9	70.8	2	2.2	65.7

C.Analysis のメニュー [分析→多変量解析他→判別手法→判別分析] を選択すると以下の分析実行メニューが表示されます。

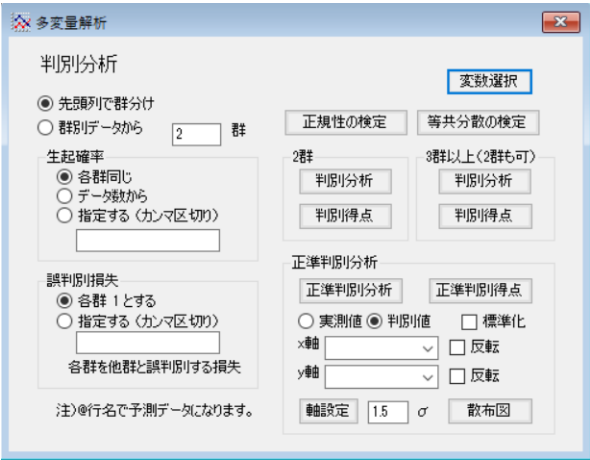


図2 判別分析実行画面

判別分析.txt (p1) は以下の通りです。ここで、「合否」の合格は 1、不合格は 2 にしています。

データ編集 判別分析.txt				
	合否	勉強時間	平均点	
▶ 1	1	5.6	70.2	
2	1	5.9	74.2	
3	1	4.1	72.7	
4	1	5.1	84.9	
5	1	5.0	93.0	
6	1	3.2	80.5	
7	1	4.3	62.7	
8	1	4.8	85.4	
9	1	3.3	84.3	
10	1	5.3	64.8	
11	1	5.9	60.7	

図3 判別分析.txt (p1)

「変数選択」は先頭に「合否」を指定して、「先頭列で群分け」にすればよいので、「All」ですべてを選びます。群は合否の 2 群なので、2 群のところの「判別分析」ボタンをクリックすると、以下の結果が表示されます。

	判別関数	標準化係数	F検定値	自由度	確率値
▶ 勉強時間	2.246	2.621	19.882	1,27	0.0001
平均点	0.201	2.279	15.027	1,27	0.0006
定数項	-23.019	-0.379			
マハラノビスの距離	5.682				
誤判別確率	1群を2群と	2群を1群と			
理論から	0.117	0.117			
実測から	0.077	0.059			
判別関数・確率	1群予測	2群予測	1群予測	2群予測	
1群実測	12	1	0.923	0.077	
2群実測	1	16	0.059	0.941	

図4 分析結果

判別関数の係数は、分析結果の「判別関数」のところです。結果への影響の強さは、隣の

「標準化係数」で分かります。重回帰分析のところでも述べたように、標準化係数は絶対値が大きいほど重要です。ここではほぼ同じですが、少し勉強時間の影響が大きいようです。また、これらの変数の係数（影響）が0かどうか判定できます。結果の右端の確率の値が有意水準0.05より小さい場合、その変数の係数は0と異なると判断され、判定に影響がある変数になります。ここでは、共に有意水準（0.05）より小さいので必要な変数です。

次に判定の正確性ですが、これは「実測値から」の「誤判別確率」がよく使われます。ただ判別関数の係数は、その上にある「理論から」の値が小さくなるようにして決められています。「マハラノビスの距離」というのが2つの群の分布の中心間の距離を表しており、これが大きくなるほど、誤判別確率の理論値が小さくなります。その概略は以下の表で与えられます。

表1 マハラノビス距離と誤判別確率

マハラノビス距離	1	4	9	16	25
誤判別確率	0.309	0.159	0.067	0.023	0.006

理論は少し横において、我々は実測値からの値に注目して行きましょう。例えば、合格（1群）を不合格（2群）と誤判別する確率は  $1/13=0.077$ 、その逆の確率は  $1/17=0.059$  となります。これは判別得点の結果を見ても分かります。分析実行画面の2群の「判別得点」をクリックすると以下の結果になります。



	所属群	判別得点	判別群
1	1	3.651	1
2	1	5.128	1
3	1	0.784	1
4	1	5.479	1
5	1	6.880	1
6	1	0.328	1
7	1	-0.774	2
8	1	4.905	1

図5 判別得点

元々の所属群と判別得点と判別群が表示されています。この中で誤判別は、所属群が1群なのに判別得点が負の値となり、2群と判別された7番です。

これは少し進んだ話になりますが、判別分析では、予め2つの群に差を付けて判別する方法もあります。実行画面の「判別の生起確率」は元々の判別群のデータ数が大きく異なる場合に利用します。例えば、非常に難しい試験の場合、合格群は不合格群に比べて小さくなります。つまり、合格の確率は不合格の確率に比べて小さくなります。このような場合、データ数に応じて、判別に補正を加えます。また、誤判別した場合の損失（ダメージ）が大きく異なるような場合も、判別に補正を加えることが考えられます。これは実行画面の「誤判別損失」のところで設定することができます。これらについてはここまでにしておきます。

今までは2群の場合の話をしてきましたが、3群以上（2群も含む）の場合には、各群に対応した判別関数を作ります。即ち、3群の場合は3種類の判別関数を作ります。群の判別は、これらの判別関数の中にデータを代入して、最大の値となった群に属するものと判断しま

す。特に2群の場合は、2つの判別関数の大きい方の群に属するものとしますが、これまでの2群の場合の判定は、2つの判別関数の差を取って、値が正になるか負になるかで大きい方の群を判別していたわけです。3群の判別を、昔フィッシャーが使ったあやめの分類の話为例に示しておきましょう。判別分析.txt (p3) を見て下さい。

	群	がくの長さ	がくの幅	花弁の長さ	花弁の幅
1	1	5.1	3.5	1.4	0.2
2	1	4.9	3	1.4	0.2
3	1	4.7	3.2	1.3	0.2
4	1	4.6	3.1	1.5	0.2
5	1	5	3.6	1.4	0.2
6	1	5.4	3.9	1.7	0.4
7	1	4.6	3.4	1.4	0.3
8	1	5	3.4	1.5	0.2
9	1	4.4	2.9	1.4	0.2

図6 判別分析.txt (p3)

このデータの群は3つです。分析実行画面の結果は以下となります。

	1群判別関数	2群判別関数	3群判別関数	1群標準化	2群標準化	3群標準化	偏え	確率値
がくの長さ	23.544	15.698	12.446	19.496	12.999	10.306	0.938	0.0103
がくの幅	23.588	7.073	3.685	10.281	3.083	1.606	0.766	0.0000
花弁の長さ	-16.431	5.211	12.767	-29.005	9.200	22.537	0.669	0.0000
花弁の幅	-17.398	6.434	21.079	-13.262	4.904	16.067	0.743	0.0000
定数項	-85.210	-71.754	-103.270	41.870	68.900	53.980		
群間の分離	Wilks's λ	0.023	確率値	0.0000				
マハラノビスの距離	1群	2群	3群					
1群	0.000	89.864	179.385					
2群	89.864	0.000	17.201					
3群	179.385	17.201	0.000					
誤判別確率	1群を他群と	2群を他群と	3群を他群と					
実測から	0.000	0.040	0.020					
判別関数・確率	1群予測	2群予測	3群予測	1群予測	2群予測	3群予測		
1群実測	50	0	0	1.000	0.000	0.000		
2群実測	0	48	2	0.000	0.960	0.040		
3群実測	0	1	49	0.000	0.020	0.980		

図7 3群の分析結果

「判別得点」をクリックすると結果は以下の通りです。

	所属群	1群	2群	3群	判別群
132	3	50.470	125.359	132.921	3
133	3	1.231	91.857	104.569	3
134	3	19.271	83.177	82.186	2
135	3	3.369	80.586	83.235	3
136	3	26.601	116.928	129.977	3
137	3	9.549	95.818	109.752	3
138	3	16.910	90.884	95.966	3

図8 3群の判別得点

間違った134番の判別関数では2群が一番大きくなっています。

最後に判別分析についてまとめておきましょう。

## 判別分析まとめ

判別分析の目的

2 群（多群）を判別する最適な 1 次式を求める。

$z = b_1 \text{ 勉強時間} + b_2 \text{ 平均点} + b_0$  判別関数

判別関数の係数は？ → 判別関数の欄

判別関数で群を分けるのは？ → 判別の分点 0（多群の場合値が最大の群）

判定に影響を与える変数は？ → 標準化係数の絶対値の大きい変数

各係数の有効性は？（要正規性・等共分散性）→ 確率の欄（係数が 0 と異なるかの検定）

誤判別の程度は？ → 誤判別確率（実測と理論）（理論値は要正規性・等共分散性）

マハラノビス距離とは → どの程度 2 群が離れているかを表わす指標

マハラノビス距離	1	4	9	16	25
誤判別確率	0.309	0.159	0.067	0.023	0.006

データ毎の判別関数の値と判別状況 → 判別得点

### 問題 1

判別分析.txt (p2) は、適性の有無の判定（有：1，無：2）と適性検査の結果と S P I の結果を与えたデータである。判定を適性検査と S P I で予測する判別分析を行い、以下の問いに答えよ。但し、事象の生起確率は各群同じ、誤判別損失は 2 群とも 1 とすること。

1) 判別関数を求めよ。

判別得点 = [ ] 適性検査 + [ ] S P I + [ ]

2) どちらの変数が判定に影響があると思われるか。[適性検査・S P I]

3) 実測値から求めた誤判別の確率は？

適性有りを無しと [ ] 適性無しを有りと [ ]

4) 先頭（1 番）の人の判別得点はいくらか。[ ]

5) 適性検査 50 点，S P I 55 点の人の判別得点はいくらか、またその人の適性の有無を判定せよ。判別得点 [ ] 適性 [有り・無し]

### 問題 2

判別分析.txt (p3) はあやめの種類をがくの長さ、花弁の長さ、花弁の幅で 3 群に分類したデータである。あやめの群を他の変数の 1 次式で判別する 3 群以上の判別分析を行い、以下の問題に答えよ。但し、設定は前問と同じとする。

1) 3 つの判別得点の式を求めよ。

判別得点 1 = [ ] がくの長さ + [ ] がくの幅  
+ [ ] 花弁の長さ + [ ] 花弁の幅 + [ ]

判別得点 2 = [ ] がくの長さ + [ ] がくの幅  
+ [ ] 花弁の長さ + [ ] 花弁の幅 + [ ]

判別得点 3 = [ ] がくの長さ + [ ] がくの幅  
+ [ ] 花弁の長さ + [ ] 花弁の幅 + [ ]

2) 実測値から求めた誤判別確率はいくらか。

群1を他と [                    ] 群2を他と [                    ] 群3を他と [                    ]

3) 先頭のデータの3つの判別得点を求めよ。

判別得点1 [                    ] 判別得点2 [                    ] 判別得点3 [                    ]

これは何群に予想されたか。 [1群・2群・3群]

#### 問題1 解答 (判別分析.txt (p2))

1) 判別関数を求めよ。

判別得点 = [ -0.190 ] 適性検査 + [ 0.645 ] S P I + [ -20.467 ]

2) どちらの変数が判定に影響があると思われるか。 [適性検査・S P I]

3) 実測値から求めた誤判別の確率は？

適性有りを無しと [ 0.053 ] 適性無しを有りと [ 0.095 ]

4) 先頭 (1番) の人の判別得点はいくらか。 [ -5.118 ]

5) 適性検査 50点, S P I 55点の人の判別得点はいくらか、またその人の適性の有無を判定せよ。 判別得点 [ 5.508 ] 適性 [有り]・無し]

#### 問題2 解答 (判別分析.txt (p3))

1) 3つの判別得点の式を求めよ。

判別得点1 = [ 23.544 ] がくの長さ + [ 23.588 ] がくの幅  
+ [ -16.431 ] 花卉の長さ + [ -17.398 ] 花卉の幅 + [ -85.210 ]

判別得点2 = [ 15.698 ] がくの長さ + [ 7.073 ] がくの幅  
+ [ 5.211 ] 花卉の長さ + [ 6.434 ] 花卉の幅 + [ -71.754 ]

判別得点3 = [ 12.446 ] がくの長さ + [ 3.685 ] がくの幅  
+ [ 12.767 ] 花卉の長さ + [ 21.079 ] 花卉の幅 + [ -103.270 ]

2) 実測値から求めた誤判別確率はいくらか。

群1を他と [ 0 ] 群2を他と [ 0.040 ] 群3を他と [ 0.020 ]

3) 先頭のデータの3つの判別得点を求めよ。

判別得点1 [ 90.940 ] 判別得点2 [ 41.644 ] 判別得点3 [ -4.808 ]

これは何群に予想されたか。 [1群]・2群・3群]