

Linear and Quadratic Regression

Probabilistic Approach

林哲緯

February 2, 2025

監督式學習中的群組分析，有別於輸入與輸出之間存在確定性關係，輸出是由模型參數完全決定的 (Deterministic)，這次從群組界線的後驗機率相等的角度切入，建立群組界線以劃分資料空間。此方法假設輸入與輸出之間存在隨機性，以機率分佈的形式表達，即是機率性方式 (Probabilistic Approach)，包括線性判別分析 (Linear Discriminant Analysis, LDA)、二次判別分析 (Quadratic Discriminant Analysis, QDA) 與 K-鄰近法 (K Nearest Neighbors, KNN)。

1 線性判別分析 (LDA)

假設 \mathbf{X} 代表多變量資料樣本變數， G 代表群組的類別變數，則後驗機率表示為 $P(G | \mathbf{X})$ ，也就是說當資料在給定 $\mathbf{X} = \mathbf{x}$ 的條件下，該資料屬於群組 $G = k$ 的機率。分別計算資料屬於不同群組的機率後，得到最大機率的群組，以此方式來推斷資料應該屬於哪個群組。因此需要計算最大後驗機率來作為群組判別的依據，如式 (1)，也稱之為判別式分析 (Discriminant Analysis)。

$$G(\mathbf{x}) = \arg \max_k \log Pr(G = k | \mathbf{X} = \mathbf{x}) \quad (1)$$

由於後驗機率 $P(G | \mathbf{X})$ 較難計算，因此須透過貝氏定理的幫助，可以得出

$$P(G = k | \mathbf{X}) = \frac{P(\mathbf{X} | G = k)P(G = k)}{\sum_l P(\mathbf{X} | G = l)P(G = l)} \quad (2)$$

其中 $P(\mathbf{X} | G = k)$ 表示第 k 組資料發生的機率密度函數，而 $P(G = k)$ 代表群組 k 發生的機率。

本文假設群組的機率密度函數 $P(X | G = k) = f_k(\mathbf{X})$ 服從多變量常態分配，寫成

$$f_k(\mathbf{X}) = \frac{1}{(2\pi)^{p/2} |\sum_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{X}-\boldsymbol{\mu}_k)^T \sum_k^{-1} (\mathbf{X}-\boldsymbol{\mu}_k)} \quad (3)$$

其中資料變數 $\mathbf{X} \in R^p$ ， μ_k 代表 \sum_k 分別代表第 k 群資料常態假設的均值與共變異矩陣。為了簡化問題，假設所有群組的共變異矩陣都相等，即 $\sum_k = \sum, \forall k$ 。加入貝氏定理與資料的常態假設後，可以將式 (1) 改寫成

$$\begin{aligned} G(\mathbf{x}) &= \arg \max_k \log Pr(G = k | \mathbf{X} = \mathbf{x}) \\ &= \arg \max_k \log (Pr(\mathbf{X} = \mathbf{x} | G = k) Pr(G = k)) \\ &= \arg \max_k \mathbf{x}^T \sum^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \sum^{-1} \boldsymbol{\mu}_k + \log Pr(G = k) \end{aligned}$$

第一行為一筆新資料 $\mathbf{X} = \mathbf{x}$ 來自哪一個群組的機率為最高，經過貝氏定理 (2) 的轉換並去除與組別 k 無關的分母，即為算式第二行。再將式 (3) 假設的常態函數代入 (共變異矩陣相同)，同樣去除與組別 k 無關的項目，即為算式第三行，其中的目標函數又稱為線性判別式函數 (Linear Discriminant Function)，即

$$\delta_k(\mathbf{x}) = \mathbf{x}^T \sum^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^T \sum^{-1} \boldsymbol{\mu}_k + \log Pr(G = k) \quad (4)$$

在式 (4) 中，除了資料 \mathbf{x} 外，其餘皆為未知，但我們仍可利用已知的資料估計這些值。

$$\hat{\mu}_k = \sum_{G \in k} \frac{(x)_i}{N_k} \quad (5)$$

μ_k 的估計為在第 k 個群組中，觀測到樣本資料的樣本平均。其中， N_k 為第 k 個群組的資料總數。

$$\hat{\Sigma}_k = \sum_{k=1}^K \sum_{G \in k} \frac{(\mathbf{x}_i - \hat{\mu}_k)(\mathbf{x}_i - \hat{\mu}_k)^T}{N - K} \quad (6)$$

Σ_k 的估計為各組資料算出來的樣本共變異矩陣的加權平均。其中， N 代表所有群組的資料總數， K 代表群組的個數。

$$Pr(G = k) \approx \hat{\pi}_k = \frac{N_k}{N} \quad (7)$$

$Pr(G = k)$ 的估計為第 k 個群組的資料總數佔所有群組資料總數的比例。

監督式學習的群組分析，其幾何意義為在資料所在的 \mathbb{R}^p 空間切割出 K 個領域 (K 是群組數)，切割的依據當然是給定的 N 筆已知資料及其群組別。而判別新資料的群組別時，則是依據資料集落在哪一個群組類別。

在繪製分界線時，必須先找出兩群組間的共同條件。因此，利用機率相等的概念，來建立這條分界線，以 k 、 j 兩群體為例，若我們觀測到的新資料落在分界線上，則它屬於 k 群體的機率會等於它屬於 j 群體的機率，也就是說，我們無法明確的判定到底該筆資料是屬於哪個群體。數學式的表達如式 (8)。

$$Pr(G = k|X = \mathbf{x}) = Pr(G = j|X = \mathbf{x}) \quad (8)$$

在後驗機率相等的群組分界原則下，結合由貝式定理 (2) 及資料的常態分配假設 (3)，分界線的函數可以從以下的轉換得到：

$$\begin{aligned} \log \frac{Pr(G = k|X = \mathbf{x})}{Pr(G = j|X = \mathbf{x})} &= \log \frac{f_k(\mathbf{x})}{f_j(\mathbf{x})} + \log \frac{Pr(G = k)}{Pr(G = j)} \\ &= \log \frac{Pr(G = k)}{Pr(G = j)} - \frac{1}{2} (\mu_k + \mu_j)^T \Sigma^{-1} (\mu_k - \mu_j) + \mathbf{x}^T \Sigma^{-1} (\mu_k - \mu_j) \\ &= 0 \end{aligned} \quad (9)$$

從資料的後驗機率相等，得到式 (9) 的線性方程式，變可以用來判斷資料的群組劃分。

2 資料分群

本節將自行生成兩個至三個群組的資料，其中，群組的資料皆服從於二元常態分配。依據不同的資料特性，例如：距離遠近、共變異矩陣以及樣本數大小，分別利用 LDA、QDA 進行群組的判別。並且，將原始資料分成訓練資料及測試資料，觀察它們的誤判率，步驟如下：

- I. 隨機選取原始資料的 80 % 為訓練資料，20 % 為測試資料。
- II. 以訓練資料進行參數的估計，並將測試資料代入，觀察預測的誤判率。
- III. 重複 100 次步驟 I 及步驟 II，並取平均。

3 兩個群組

以下將針對兩個群組的資料，利用不同的方法，對其進行檢視及群組的判別，共有七種不同型態的資料。其中，令 Group A = 0、Group B = 1。

DATA 1.

Group A :

$$n_1 = 200, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數及共變異矩陣相同，但兩組間的距離較遠。並且，Group A 及 Group B 中的樣本皆互相獨立。由圖 1 可以看出兩群組的部分資料點有些許的重疊，分界較明顯，且兩群資料皆無明顯的趨勢。

LDA & QDA

圖 1 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 1，LDA 及 QDA 的分界線僅有些微差距。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 6.29 %，測試資料的誤判率為 6.70 %；QDA 訓練資料的誤判率為 4.73 %，測試資料的誤判率為 5.36 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，雖然只有些微的差距，但 QDA 的誤判率都些微低於 LDA。因此，推測在這種資料形態下，QDA 的判別能力是比 LDA 好的。

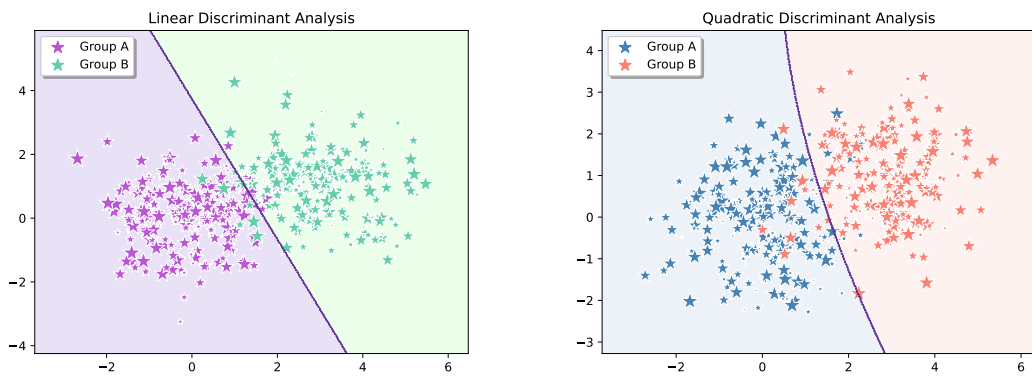


Figure 1: 原始資料的群組分界線 [DATA 1]

DATA 2.

Group A :

$$n_1 = 200, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0.7 \\ 0.7 & 1 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數相同，但兩組間的距離較遠且共變異數矩陣不同。並且，Group A 的兩變數皆互相獨立，Group B 的樣本則為相依。由圖 2 可以看出兩群組的資料點只有些許的重疊，分界較明顯，並且 Group B 有明顯的趨勢。

LDA & QDA

圖 2 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 2，可以看到 LDA 及 QDA 的分界線有明顯的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 9.38 %，測試資料的誤判率為 9.57 %；QDA 訓練資料的誤判率為 5.68 %，測試資料的誤判率為 5.91 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，LDA 的誤判率都比 QDA 還高出許多。因此，推測在這種資料形態下，QDA 的判別能力是比 LDA 好的。

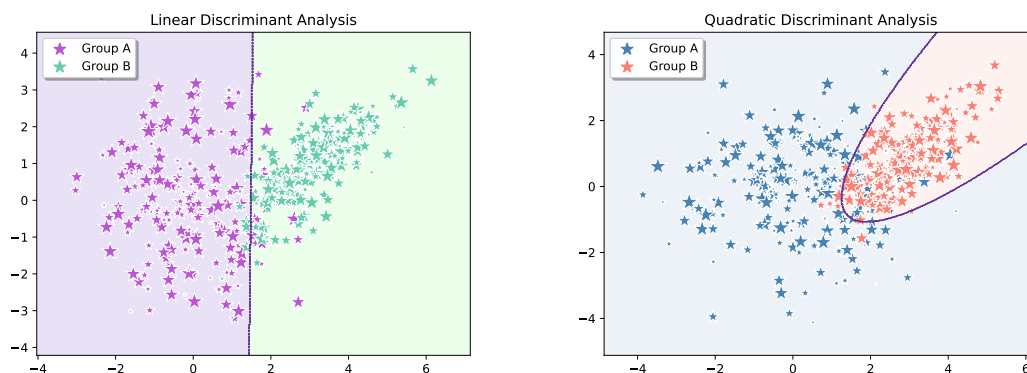


Figure 2: 原始資料的群組分界線 [DATA 2]

DATA 3.

Group A :

$$n_1 = 200, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數相同，但兩組間的距離較遠且共變異數矩陣不同。並且，Group A 及 Group B 中的變數皆互相獨立，只是 Group B 變異較大，且資料點較分散。由圖 3 可以看出兩群組的資料點只有些許的重疊，分界較明顯，並且 Group B 較 Group A 分散許多。

LDA & QDA

圖 3 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 3，可以看到 QDA 的分界線是一個圓圈，與 LDA 的分界線有明顯的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 9.85 %，測試資料的誤判率為 9.81 %；QDA 訓練資料的誤判率為 8.73 %，測試資料的誤判率為 7.71 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，LDA 的誤判率都些微高出 QDA。因此，推測在這種資料形態下，QDA 的判別能力是比 LDA 好的。

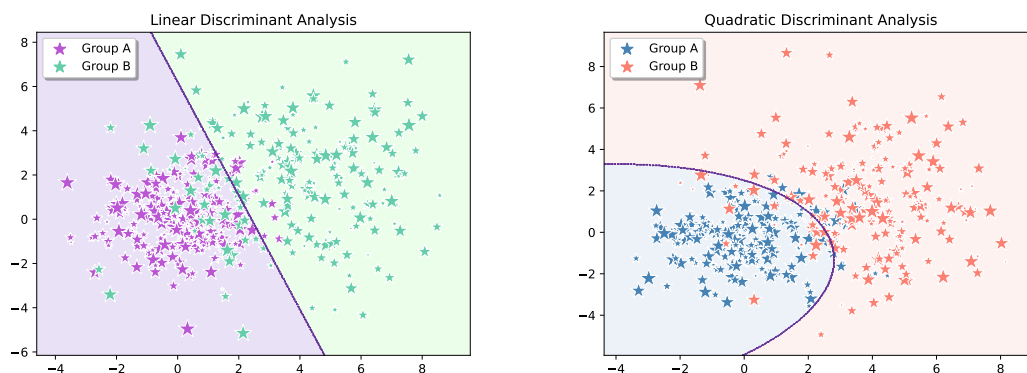


Figure 3: 原始資料的群組分界線 [DATA 3]

DATA 4.

Group A :

$$n_1 = 200, \mu_1 = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 1 & -0.7 \\ -0.7 & 1 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} 1 \\ 3 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數相同、距離較近，且兩組間的共變異數矩陣不同。並且，Group A 及 Group B 中的樣本皆為相依，Group B 資料點較集中，兩群都有明顯的趨勢。由圖 4 可以看出兩群組中間部分的資料點是重疊在一起的，呈現一個交叉的形狀。

LDA & QDA

圖 4 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 4，可以看到 QDA 與 LDA 的分界線有明顯的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 43.14 %，測試資料的誤判率為 45.66 %；QDA 訓練資料的誤判率為 18.68 %，測試資料的誤判率為 18.58 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，LDA 的誤判率都明顯比 QDA 高出許多。因此，推測在這種資料形態下，QDA 的判別能力是比 LDA 好的。

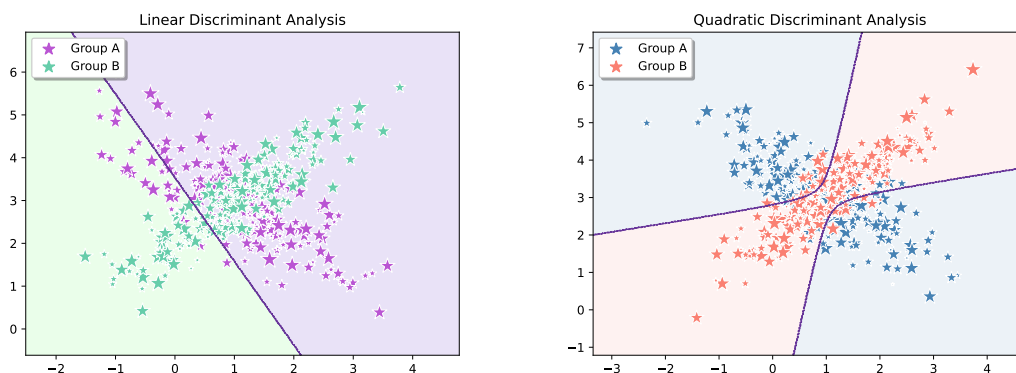


Figure 4: 原始資料的群組分界線 [DATA 4]

DATA 5.

Group A :

$$n_1 = 200, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數及共變異數矩陣皆相同，且兩群資料點分散程度相近。兩組資料皆為相互獨立。由圖 5 可以看出兩群組的資料點幾乎是完全重疊在一起的。

LDA & QDA

圖 5 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 5，可以看到 QDA 與 LDA 的分界線有些微的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 24.67 %，測試資料的誤判率為 25.45 %；QDA 訓練資料的誤判率為 25.02 %，測試資料的誤判率為 26.23 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，QDA 的誤判率都比 LDA 些微高出一點點。因此，推測在這種資料形態下，兩種判別方法可能表現皆差不多。

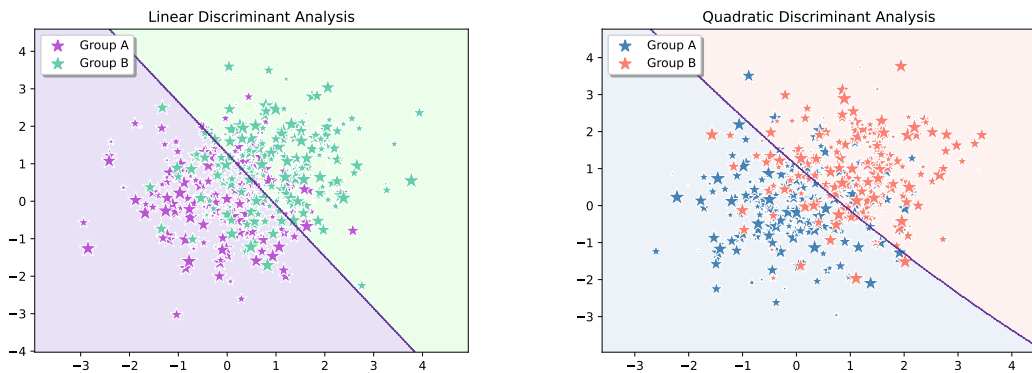


Figure 5: 原始資料的群組分界線 [DATA 5]

DATA 6.

Group A :

$$n_1 = 500, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$$

Group B :

$$n_2 = 150, \mu_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數及共變異矩陣不相同。兩組資料皆為相互獨立，其中 Group A 的變異較大，資料較分散。由圖 6 可以看出兩群組的資料點是完全重疊在一起的，並且，Group B 是位於 Group A 之中的。

LDA & QDA

圖 6 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 6，可以看到 QDA 的分界線是一個橢圓，與 LDA 的分界線有明顯的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 35.56 %，測試資料的誤判率為 36.76 %；QDA 訓練資料的誤判率為 20.62 %，測試資料的誤判率為 21.59 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，LDA 的誤判率都比 QDA 高。因此，推測在這種資料形態下，QDA 的判別能力可能是比 LDA 好的。

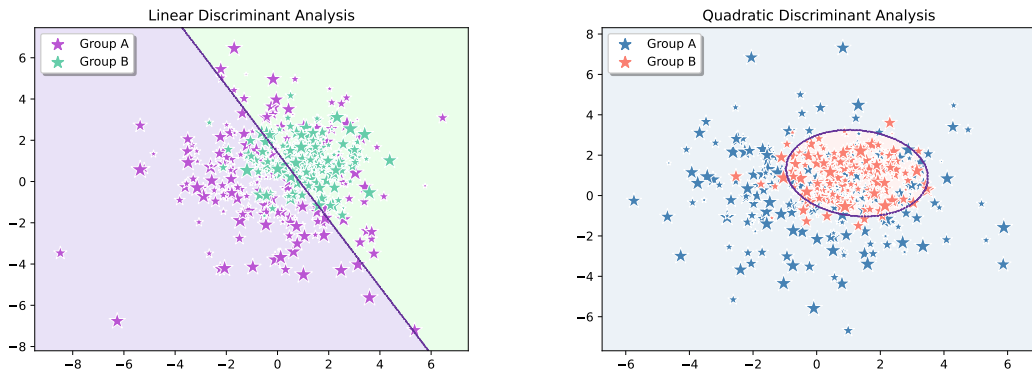


Figure 6: 原始資料的群組分界線 [DATA 6]

DATA 7.

Group A :

$$n_1 = 200, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 3 & 0.7 \\ 0.7 & 3 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} 4 \\ 2 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 3 & 0.7 \\ 0.7 & 3 \end{bmatrix}$$

此資料的型態為 Group A 及 Group B 的樣本數及共變異矩陣皆相同，距離較遠，且兩組資料皆為相依。由圖 7 可以看出兩群組的資料點是有明顯的分界，重疊的資料點較少，並且兩組資料的趨勢皆相同。

LDA & QDA

圖 7 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 7，可以看到 QDA 與 LDA 的分界線有明顯的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 10.74 %，測試資料的誤判率為 10.84 %；QDA 訓練資料的誤判率為 9.65 %，測試資料的誤判率為 10.04 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，QDA 的誤判率都比 LDA 低。因此，推測在這種資料形態下，QDA 的判別能力可能稍微比 LDA 好。

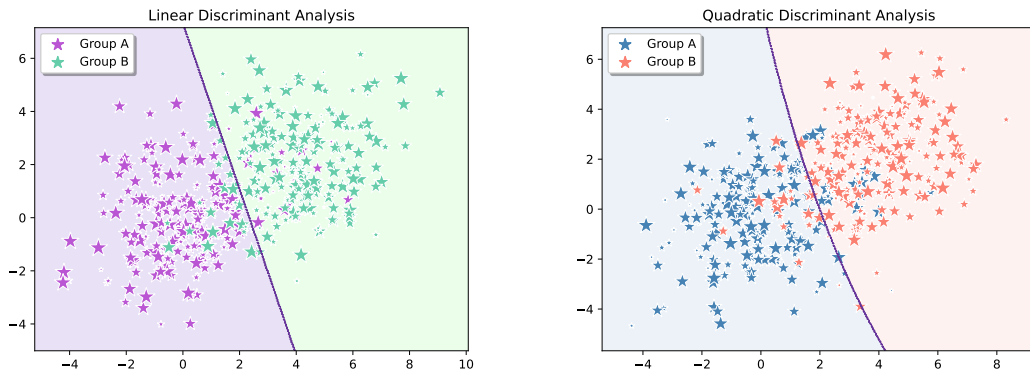


Figure 7: 原始資料的群組分界線 [DATA 7]

3.1 三個群組

以下將針對三個群組的資料，利用不同的方法，對其進行檢視及群組的判別，共有兩種不同型態的資料。其中，令 Group A = 0、Group B = 1、Group C = 2。

DATA 8.

Group A :

$$n_1 = 300, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

Group B :

$$n_2 = 200, \mu_2 = \begin{bmatrix} -4 \\ -2 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 2 & -0.7 \\ -0.7 & 1 \end{bmatrix}$$

Group C :

$$n_3 = 100, \mu_3 = \begin{bmatrix} -4 \\ 2 \end{bmatrix}, \Sigma_3 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

LDA & QDA

由圖 8 可以看出這三群組的資料點之間是有明顯的分界，重疊的資料點較少。觀察圖 8，可以看到 QDA 與 LDA 的分界線在 Group C 與 Group A 或 Group B 之間有明顯的差別。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 5.35 %，測試資料的誤判率為 5.49 %；QDA 訓練資料的誤判率為 5.33 %，測試資料的誤判率為 5.83 %。LDA 的誤判率和 QDA 誤判率皆差不多。因此，推測在這種資料形態下，兩者表現相當。

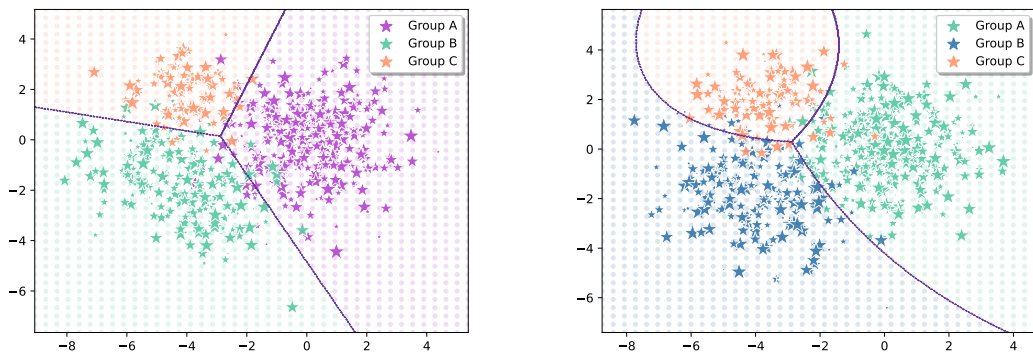


Figure 8: 原始資料的群組分界線 [DATA 8]

DATA 9.

Group A :

$$n_1 = 500, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 5 & 0.5 \\ 0.5 & 5 \end{bmatrix}$$

Group B :

$$n_2 = 300, \mu_2 = \begin{bmatrix} -5 \\ 0 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 3 & -0.2 \\ -0.2 & 3 \end{bmatrix}$$

Group C :

$$n_3 = 250, \mu_3 = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, \Sigma_3 = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$$

由圖 9 可以看出這三群組的資料點之間是有明顯的分界，重疊的資料點雖比 DATA 8 多，但重疊的情況仍算少。

LDA & QDA

圖 9 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 9，可以看到 QDA 與 LDA 的分界線在 Group A 與 Group B 之間有明顯的差別，Group A 與 Group C 之間的差異則比較不明顯。LDA 訓練資料的誤判率為 16.67 %，測試資料的誤判率為 17.21 %；QDA 訓練資料的誤判率為 15.89 %，測試資料的誤判率為 16.41 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，LDA 的誤判率都比 Q 稍微高一些。因此，推測在這種資料形態下，QDA 的判別能力可能是比 LDA 好的。

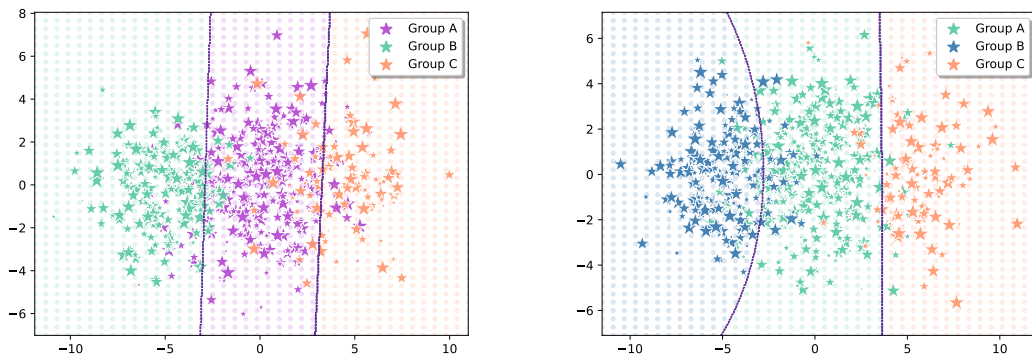


Figure 9: 原始資料的群組分界線 [DATA 9]

DATA 10.

Group A :

$$n_1 = 500, \mu_1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma_1 = \begin{bmatrix} 3 & -0.7 \\ -0.7 & 3 \end{bmatrix}$$

Group B :

$$n_2 = 300, \mu_2 = \begin{bmatrix} -5 \\ 0 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 3 & -0.7 \\ -0.7 & 3 \end{bmatrix}$$

Group C :

$$n_3 = 250, \mu_3 = \begin{bmatrix} 5 \\ 0 \end{bmatrix}, \Sigma_3 = \begin{bmatrix} 3 & -0.7 \\ -0.7 & 3 \end{bmatrix}$$

此資料的型態為 Group A、Group B 及 Group C 的共變異矩陣皆相同，距離較遠，且三組資料皆為相依。由圖 10 可以看出這三群組的資料點之間是有明顯的分界，重疊的資料點較少。

LDA & QDA

圖 10 為原始資料利用 LDA、QDA 所繪製的分界線。觀察圖 10，可以看到 QDA 在 Group A 與 Group B 及 Group A 與 Group C 之間的分界線較 LDA 彎曲一些。利用訓練資料及測試資料的誤判率來比較，LDA 訓練資料的誤判率為 8.90 %，測試資料的誤判率為 8.97 %；QDA 訓練資料的誤判率為 9.44 %，測試資料的誤判率為 9.80 %。從誤判率可以發現，不論是在訓練資料，亦或是測試資料，QDA 的誤判率都比 LDA 高一些。因此，推測在這種資料形態下，LDA 的判別能力可能是比 QDA 好的。

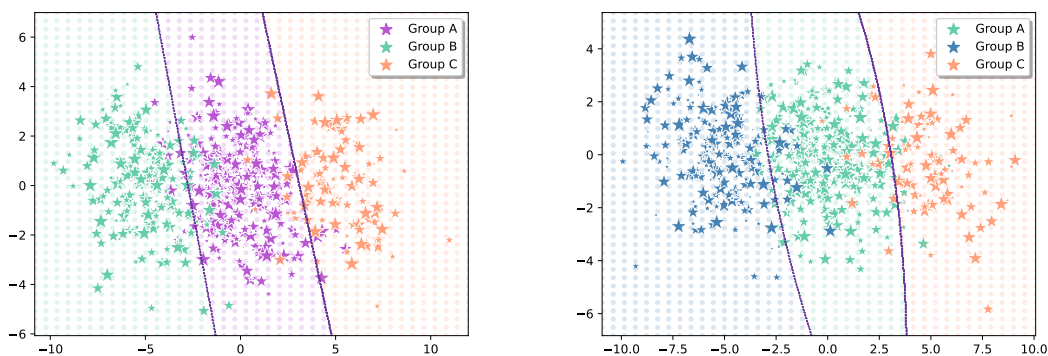


Figure 10: 原始資料的群組分界線 [DATA 10]

4 方法比較

本節將整理以上兩種機率性方法，在不同資料下的誤判率，並推測出針對不同資料，較適用於何種方法來進行判別預測。

Table 1: 兩群資料的機率性方法的誤判率比較

訓練資料誤判率 (%)	DATA 1	DATA 2	DATA 3	DATA 4	DATA 5	DATA 6	DATA 7
線性判別分析 (LDA)	6.29	9.38	9.85	43.14	24.67	35.56	10.74
二次判別分析 (QDA)	4.73	5.68	8.73	18.68	25.02	20.62	9.65
測試資料誤判率 (%)	DATA 1	DATA 2	DATA 3	DATA 4	DATA 5	DATA 6	DATA 7
線性判別分析 (LDA)	6.70	9.57	9.81	45.66	25.45	36.76	10.84
二次判別分析 (QDA)	5.36	5.91	7.71	18.58	26.23	21.59	10.04

Table 2: 兩個群組的資料型態

Group		DATA 1	DATA 2	DATA 3	DATA 4	DATA 5	DATA 6	DATA 7
A	樣本數	200	200	200	200	200	500	200
	變異數	1	2	1	1	1	5	3
	變數間的關係	獨立	獨立	獨立	相依	獨立	獨立	相依
B	樣本數	200	200	200	200	200	150	200
	變異數	1	1	5	1	1	1	3
	變數間的關係	獨立	相依	獨立	相依	獨立	獨立	相依
兩母體的距離		較遠	較遠	較遠	較近	較近	較近	較遠
共變異矩陣		相同	不同	不同	不同	相同	不同	相同