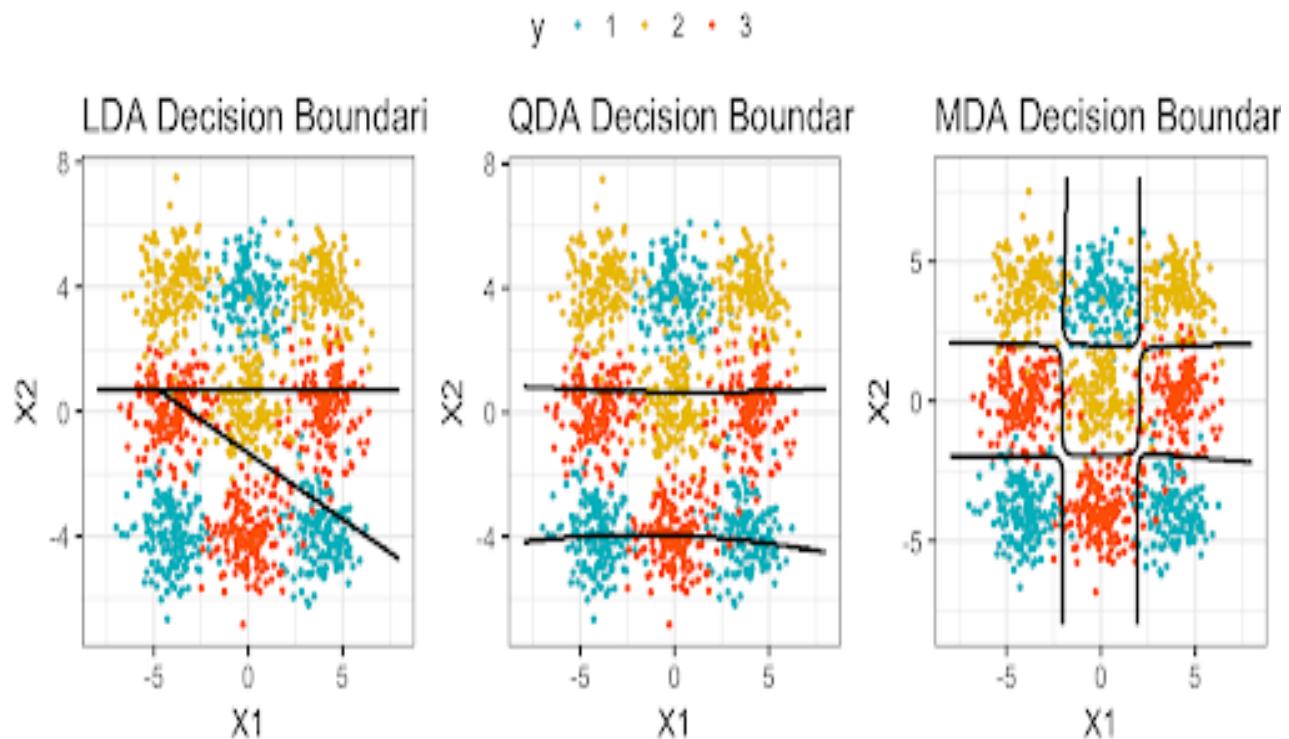


# Quadratic Discriminant Analysis

Ridwan Olawin

---



## Abstract

Classification in machine learning is a supervised learning approach in which a computer program learns from the data given to it and makes new observations or classifications. There are a variety of classification algorithms out there but the one factor that stands out amongst them is their individual accuracy and efficiency. This has led to scientists making assumptions about data before implementing the algorithm to get the best classifier possible. One of those methods led to the creation of generative classifiers which uses data to make further assumptions about how to classify itself. In this paper, we explore in detail Quadratic Discriminant Analysis (QDA), a generative classifier, its advantages, and disadvantages over similar algorithms to itself.

## Background

Generative classifiers are typically used to generate models of joint probability distribution. This means there exists an input and target variable (X and y). The solution would be a distribution that could generate new input variables with their respective targets. The general name of the model to be extended in this work is called a Gaussian Discriminant Analysis (GDA) model. For a model of this category, a major assumption here is that the class conditional densities are normally distributed.

$$P(\mathbf{x} | t = c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c),$$

Where  $\boldsymbol{\mu}$  is the class-specific mean vector and  $\boldsymbol{\Sigma}$  is the class-specific covariance matrix. Using bayes classifier, we can calculate the class posterior:

$$\overbrace{P(t = c | \mathbf{x}, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}^{\text{class posterior}} = \frac{\overbrace{P(\mathbf{x} | t = c, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)}^{\text{class-conditional density}} \overbrace{P(t = c)}^{\text{class prior}}}{\sum_{k=1}^K P(\mathbf{x} | t = k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) P(t = k)}.$$

*Equation 1*

We can then classify  $\mathbf{x}$  into class

$$\hat{h}(\mathbf{x}) = \operatorname{argmax}_c P(t = c | \mathbf{x}, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c).$$

*Equation 2*

Given the above conditions, we are typically trying to reduce the error rate:

$$\text{Error rate} = L(h) = P(h(\mathbf{x}) \neq y)$$

- $y$  is the actual and  $h(\mathbf{x})$  is the predicted value of our target

The bayes classifier depends on unknown quantities so we need to use data to find some approximation to the rule. This resulted in a concept known as **Gaussian Maximum Likelihood Classification (GMLC)**.

# Related Work

## Gaussian Maximum Likelihood Classification

This method assumes that each class (target) has a gaussian distribution and the estimates the distribution from the data, then classifies each new observation to the class with maximum likelihood. There are two main methods that fall under this classification known as Linear Discriminant Analysis (LDA) and QDA. It is important to understand in detail the background of LDA before exploring QDA.

# Methodology

## Linear Discriminant Analysis

LDA can be used to perform supervised dimensionality reduction (as learned in class). In this case however, it is used to model the class conditional distribution of data  $P(x|y=k)$  for each class k.

$$P(x|y=k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_k)^t \Sigma_k^{-1} (x - \mu_k)\right)$$

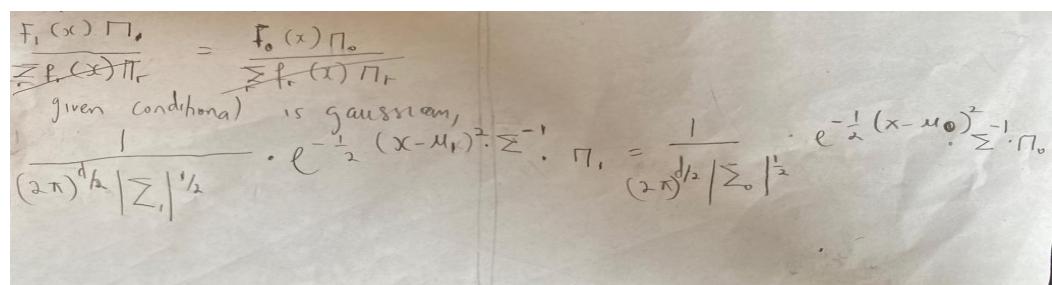
Given this conditional, we simply want a way to draw the decision boundary between our classes:

$$\{ x \mid P(y=1|x=x) == P(y=0|x=x) \}$$

We are simply trying to solve for

$$P(y=1|x=x) == P(y=0|x=x)$$

**A major assumption in LDA is that the class covariances are the same.** Figures 1 and 2 show the derivation of the decision boundary.



The image shows a handwritten derivation of the LDA decision boundary formula. It starts with the ratio of two Gaussian probability density functions:

$$\frac{f_1(x) \pi_1}{f_0(x) \pi_0} = \frac{f_0(x) \pi_0}{f_1(x) \pi_1}$$

Given the condition that the distributions are Gaussian, the formula is simplified to:

$$\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \cdot e^{-\frac{1}{2} (x - \mu_1)^t \Sigma^{-1} (x - \mu_1)} \cdot \pi_1 = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \cdot e^{-\frac{1}{2} (x - \mu_0)^t \Sigma^{-1} (x - \mu_0)} \cdot \pi_0$$

Figure 1

- Last assumption being made + is that

$$\Sigma = \Sigma_0 = \Sigma_1 \quad (\text{covariance of both classes are the same})$$

after this assumption, denominators are the same -

$$e^{-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)} = e^{-\frac{1}{2}(x - \mu_2)^T \Sigma^{-1} (x - \mu_2)}$$

• Take log of both sides, and moving RHS to the LHS

$$-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \log(\pi_1) + \frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0) - \log(\pi_0) = 0$$

$$-\frac{1}{2}x^T \Sigma^{-1} x - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + x^T \Sigma^{-1} \mu_1 + \frac{1}{2}x^T \Sigma^{-1} x + \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - x^T \Sigma^{-1} \mu_0 + \log\left(\frac{\pi_1}{\pi_0}\right) = 0$$

$$x^T (\Sigma^{-1} \mu_0 - \Sigma^{-1} \mu_1) + \frac{1}{2} (\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) + \log\left(\frac{\pi_1}{\pi_0}\right) = 0$$

This is just a line      This is constant

Figure 2

From the derivation, we can see the equation ending in the form:

$x^T B + a = 0$ ; which simply shows the form  $mx + c$ , the equation of a line.

The resulting graph looks of the form:

► **Linear: Decision boundaries are linear**

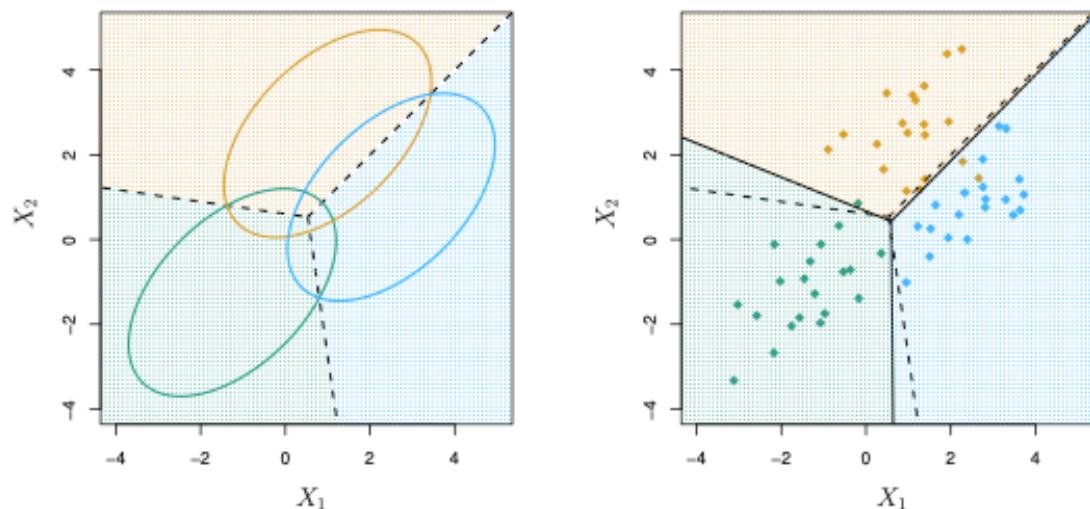


Figure 3

Taking the log of the last equation in Figure 2 would result in the decision boundary equation below:

$$\Sigma_0 = \Sigma_1 = \Sigma_2 = \dots = \Sigma_{k-1}$$

(e.g. LDA), then

$$\delta_k(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^\top \Sigma^{-1} \boldsymbol{\mu}_k + \log(\pi_k)$$

Some notations to note:

If we wish to use LDA we must calculate a common covariance, so we average all the covariances e.g.

$$\Sigma = \frac{\sum_{r=1}^k (n_r \Sigma_r)}{\sum_{r=1}^k n_r}$$

Where:

$n_r$  is the number of data points in class  $r$

$\Sigma_r$  is the covariance of class  $r$

~~$n$  is the total number of data points, and  $k$  is the number of classes~~

In practice we don't know the parameters of Gaussian and will need to estimate them using our training data.

$$\hat{\pi}_k = \hat{P}r(y = k) = \frac{n_k}{n}$$

where  $n_k$  is the number of class  $k$  observations.

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$$

$$\hat{\Sigma}_k = \frac{1}{n_k - k} \sum_{i:y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^\top$$

## Quadratic Discriminant Analysis

This method is simply an extension of LDA except one of the assumptions is being changed. The assumption is that the covariances of all classes are different. Its resulting graph is quadratic hence, its name QDA.

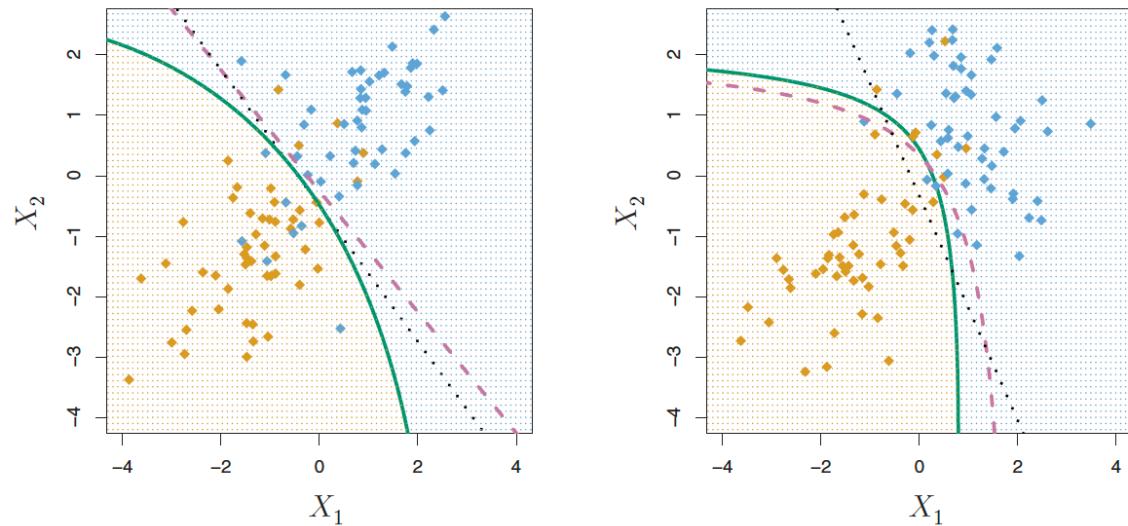


Figure 4

Figures 5 and 6 below simply expand on the different cases of what the covariance matrices could be.

For QDA we need to calculate:

$$\delta_k(\mathbf{x}) = -\frac{1}{2} \log(|\Sigma_k|) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) + \log(\pi_k)$$

Equation 3

Let's expand on the different cases of the covariant matrices

Case 1:

$$\delta_i(x) = \frac{1}{2} \log(\Pi_i) - \frac{1}{2} (x - \mu_i)^T (x - \mu_i) + \log(\Pi_i)$$

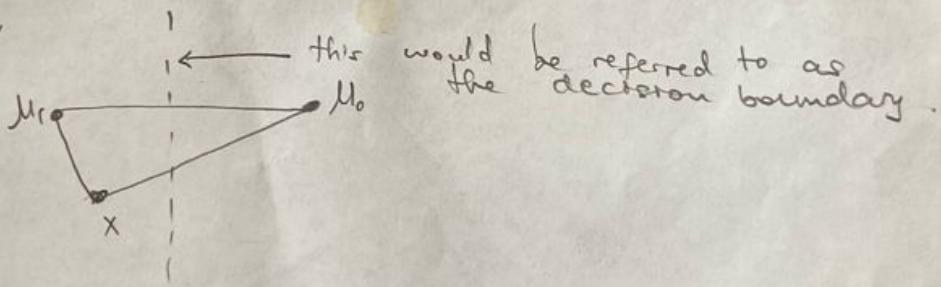
this is  
basically a  
squared  
Euclidean  
distance

another assumption,  
let's assume we  
have equal  
class distribution  
hence this  
would be  
negligible.

The conclusion here is basically is coming to understand that QDA can be looked at in the sense of squared euclidean distance from the mean of our particular data point. different classes to a

$\delta_i(x)$  Hence,

given  
 $\delta_i(x)$   
 $d_i(x)$ ,



Case 2:

$\Sigma_k \neq I$   
recall in SVD (singular-value decomposition)

$$A = U \Sigma V^T$$

$U$  is eigenvectors of  $A^T A$

$V$  " " "  $A^T A$

$\Sigma$  is a diagonal matrix such that diagonal values are eigenvalues of  $A^T A$  or  $A A^T$

Figure 5

∴ Hence in case 2 where  
 $\Sigma_k \neq I$ ,

$$\Sigma_k = U \Sigma_k^T$$

Since  $\Sigma_k$  is symmetric,

$$\Sigma_k = U \Sigma_k^T$$

recall:

$$(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \leftarrow$$
$$\Sigma_k^{-1} = (\mu \Sigma_k^T)^{-1} = \mu \Sigma^{-1} \mu^T$$

$$(x - \mu_k)^T (\mu \Sigma^{-1} \mu^T) (x - \mu_k)$$

$$(\mu^T x - \mu^T \mu_k)^T \Sigma^{-1} (\mu^T x - \mu^T \mu_k)$$

$$(\mu^T x - \mu^T \mu_k)^T \Sigma^{-\frac{1}{2}} \Sigma^{-\frac{1}{2}} (\mu^T x - \mu^T \mu_k)$$

$$(\Sigma^{-\frac{1}{2}} \mu^T x - \Sigma^{-\frac{1}{2}} \mu^T \mu_k) (\Sigma^{-\frac{1}{2}} \mu^T x - \Sigma^{-\frac{1}{2}} \mu^T \mu_k)$$

Let's think of

$\Sigma^{-\frac{1}{2}} \mu^T$  as a linear transformation.

$$x \rightarrow \Sigma^{-\frac{1}{2}} \mu^T x$$

claim here  $x$  becomes spherical  
which would reverse us back  
to case 1.

Figure 6

# Experiments and Result

To show the implementation of these classification methods, I used a Kaggle data set which contained data about passengers on board the titanic. Some exploratory analysis of the data set showed the distribution of the data including missing data sets and unclear ones.

---

	<b>Survived</b>	<b>Pclass</b>	<b>Age</b>	<b>SibSp</b>	<b>Parch</b>	<b>Fare</b>
<b>count</b>	891.000000	891.000000	714.000000	891.000000	891.000000	891.000000
<b>mean</b>	0.383838	2.308642	29.699118	0.523008	0.381594	32.204208
<b>std</b>	0.486592	0.836071	14.526497	1.102743	0.806057	49.693429
<b>min</b>	0.000000	1.000000	0.420000	0.000000	0.000000	0.000000
<b>25%</b>	0.000000	2.000000	20.125000	0.000000	0.000000	7.910400
<b>50%</b>	0.000000	3.000000	28.000000	0.000000	0.000000	14.454200
<b>75%</b>	1.000000	3.000000	38.000000	1.000000	0.000000	31.000000
<b>max</b>	1.000000	3.000000	80.000000	8.000000	6.000000	512.329200

Figure 7

My experiment process revolved around implementing both LDA and QDA from scratch using the methodology above. The goal of this was to show the improvement in accuracy QDA had over LDA in this data set. The reason why this occurred was because this data set had a relatively larger training set and higher variance. The results are below:

**Precision = 0.4257**  
**Recall = 0.5676**  
**F\_measure = 0.4865**  
**Accuracy for LDA Mine: 0.503731343283582**

Figure 8: LDA

**Precision = 0.7358**  
**Recall = 0.7027**  
**F\_measure = 0.7189**  
**Accuracy for QDA Mine: 0.7723880597014925**

Figure 9: QDA

# Conclusion

From the results above, we can see that QDA performs considerably better than LDA simply because our initial assumptions worked better for this data set. A few things to note specifically on this data set was that I noticed that being a female or a child increases the chances of survival. A higher-class ticket also improved survival compared to a third- or fourth-class ticket.

# Future Work

Another algorithm I am looking forward to exploring is Multiple Discriminant Analysis (MDA). This algorithm is a multivariate dimensionality reduction technique (Similar to LDA). However, it can also be used to support classification by yielding a compressed signal which can then be used for classification. The method also reduces the curse of dimensionality by compressing the signal (Features) down to a lower dimensional space (Similar to LDA). What makes this method unique is its focus on when three or more target variables are involved, it can compute more than one function as opposed to LDA which only computes one.

# Bibliography

Julenn. (2021, May 15). *Titanic*. Kaggle. <https://www.kaggle.com/julenn/titanic>.

YouTube. (2015, September 27). *Ali Ghodsi, Lec 2: Machine learning. classification, Linear and quadratic discriminant analysis*. YouTube. [https://www.youtube.com/watch?v=\\_m7TMkzZzus](https://www.youtube.com/watch?v=_m7TMkzZzus).

Petrik, M. (2017, February 16). *LDA, QDA, Naive Bayes*. Generative Classification Models. [https://marek.petrik.us/teaching/intro\\_ml\\_17/intro\\_ml\\_17\\_files/class5.pdf](https://marek.petrik.us/teaching/intro_ml_17/intro_ml_17_files/class5.pdf).