



Modern Machine Learning Linear Methods for Classification — Linear Discriminate Analysis

Kenneth E. Barner
Department of Electrical and Computer Engineering
University of Delaware



Problem: **Classification** is the prediction of **qualitative** responses

Consider the binary classification problem

Examples:

- Email: Spam/Not Spam
- Tumor: Benign/Malignant

Output $Y \in \{0,1\}$ $\begin{cases} 0 := \text{Negative class (e. g. Not Spam)} \\ 1 := \text{Positive class (e. g. Spam)} \end{cases}$

Objective: Divide the input space into a collection of regions labeled according to classification. **Linear methods for classification** result in **linear decision boundaries**



Assumption: Suppose there are K classes and the fitted linear model for the k^{th} indicator response variable is

$$\hat{f}_k(\mathbf{x}) = \hat{\beta}_{k0} + \hat{\boldsymbol{\beta}}_k^T \mathbf{x} \quad k = 1, 2, \dots, K$$

The decision boundary between classes k and l is the set of points for which $\hat{f}_k(\mathbf{x}) = \hat{f}_l(\mathbf{x})$

$$\{\mathbf{x}: (\hat{\beta}_{k0} - \hat{\beta}_{l0}) + (\hat{\boldsymbol{\beta}}_k^T - \hat{\boldsymbol{\beta}}_l^T) \mathbf{x} = 0\}$$

Note: The decision boundaries are a set of hyper planes \Rightarrow the input space is divided into regions with piecewise hyper plane decision boundaries

- This regression approach is within the class of methods that model **discrimination functions** for each class $\delta_k(\mathbf{x})$
 - \mathbf{x} is classified to the class with the largest value for its discriminant function
 - Posterior probability models $\Pr(G = k | \mathbf{X} = \mathbf{x})$ are in this class
 - If $\delta_k(\mathbf{x})$ or $\Pr(G = k | \mathbf{X} = \mathbf{x})$ are linear in \mathbf{x} , then the decision boundaries are linear

Notation: $G(\mathbf{x})$ is the class predictor

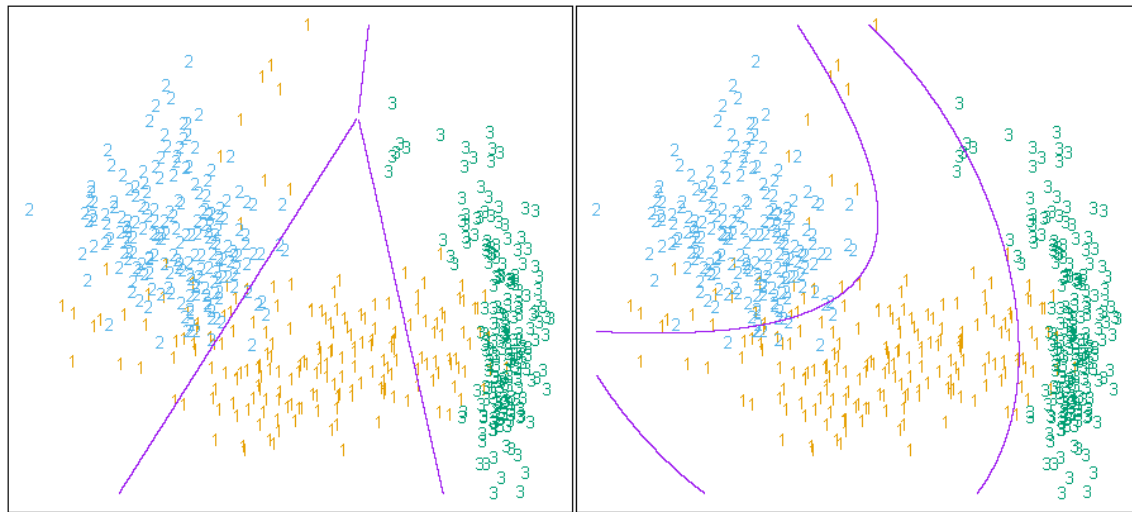


Example: A three class example is shown.

- Left Plot: linear decision boundaries obtained through linear discriminant analysis (LDA)
- Right Plot: quadratic decision boundaries obtained by expanding the observation space to five dimensions:

$$X_1, X_2, X_1X_2, X_1^2, X_2^2$$

Note: Linear boundaries in the five-dimension space are quadratic boundaries in the original space





Linear Discriminant Analysis (LDA): Methodology utilizes the class posteriors $\Pr(G|\mathbf{X})$ for optimal classification and a Gaussian model.

- Let the set of classes be $G \in \{1, \dots, K\}$
- Define $f_k(\mathbf{x})$ as the class-conditional density of \mathbf{X} in class $G = k$
 $\Rightarrow f_k(\mathbf{x}) = P(\mathbf{X} = \mathbf{x}|G = k)$, where $k \in \{1, \dots, K\}$
- Let π_k be the prior probability of class k
 $\Rightarrow \pi_k = P(G = k)$, where $k \in \{1, \dots, K\}$

Recall **Bayes theorem**

$$\begin{aligned} P(G = k|\mathbf{X} = \mathbf{x}) &= \frac{P(\mathbf{X} = \mathbf{x}|G = k)P(G = k)}{P(\mathbf{X} = \mathbf{x})} \\ &= \frac{P(\mathbf{X} = \mathbf{x}|G = k)P(G = k)}{\sum_{\ell=1}^K P(\mathbf{X} = \mathbf{x}|G = \ell) P(G = \ell)} \end{aligned}$$



Using the previous definitions

$$P(G = k | \mathbf{X} = \mathbf{x}) = \frac{f_k(\mathbf{x})\pi_k}{\sum_{\ell}^K f_{\ell}(\mathbf{x})\pi_{\ell}} \quad (*)$$

Observation: The denominator in $(*)$ is not a function of k . Thus to determine the class, we need only be concerned with maximizing the numerator.

- Need to determine an appropriate distribution model, e.g., Gaussian, mixtures of Gaussians, nonparametric or other densities

Assumption: Choose a Gaussian mixture model for $f_k(\mathbf{x})$

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}$$

Also assume that the classes have the same covariance matrix $\boldsymbol{\Sigma}_k = \boldsymbol{\Sigma} \forall k$



In practice we do not know the Gaussian distribution parameter values

⇒ estimate them using the training data

$\hat{\pi}_k = N_k/N$, where N_k and N are the number of class- k observation and the total number of observations, respectively

$\hat{\mu}_k = \sum_{g_i=k} \mathbf{x}_i / N_k$ (average of \mathbf{x} over each category)

$\hat{\Sigma} = \sum_{k=1}^K \sum_{g_i=k} (\mathbf{x}_i - \hat{\mu}_k)(\mathbf{x}_i - \hat{\mu}_k)^T / (N - K)$



We can show that LDA produces linear decision boundaries by taking the log-odds for classes k and ℓ

$$\begin{aligned}\log \frac{P(G = k|X = \mathbf{x})}{P(G = \ell|X = \mathbf{x})} &= \log \frac{f_k(x)}{f_\ell(x)} + \log \frac{\pi_k}{\pi_\ell} \\ &= \log \frac{\pi_k}{\pi_\ell} - (\boldsymbol{\mu}_k + \boldsymbol{\mu}_l)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l) \\ &\quad + \mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l)\end{aligned}$$

Let $\beta_0 = \log \frac{\pi_k}{\pi_\ell} - (\boldsymbol{\mu}_k + \boldsymbol{\mu}_l)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l)$ and $\boldsymbol{\beta} = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l)$, then

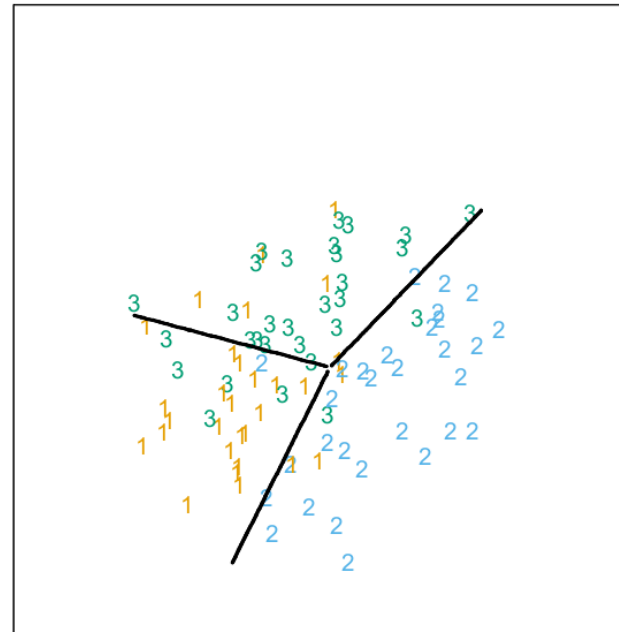
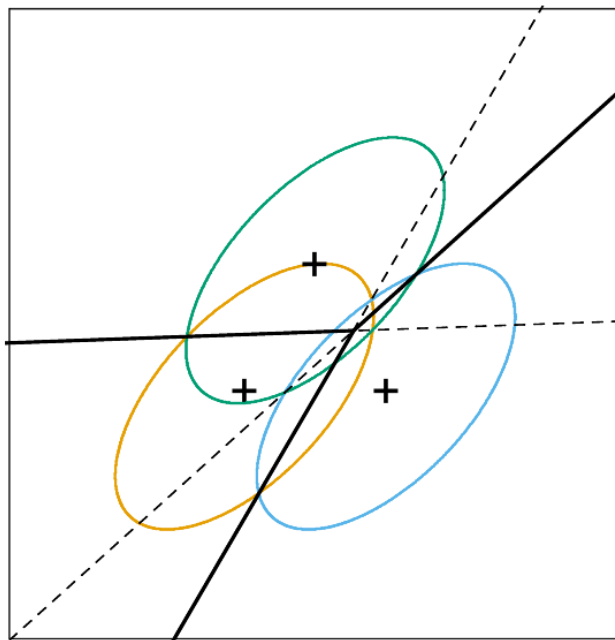
$$\log \frac{P(G = k|X = \mathbf{x})}{P(G = \ell|X = \mathbf{x})} = \beta_0 + \mathbf{x}^T \boldsymbol{\beta},$$

This is a linear equation, $\beta_0 + \mathbf{x}^T \boldsymbol{\beta} = 0$, that defines the LDA decision boundary between classes k and ℓ

**Example:**

Left — Three Gaussian distributions with the same covariance and different means. Contours enclosing 95% of the density are shown. Broken lines show class pair decision boundaries. Solid lines show boundaries separating all three classes.

Right — 30 samples from each class and the fitted LDA boundaries.





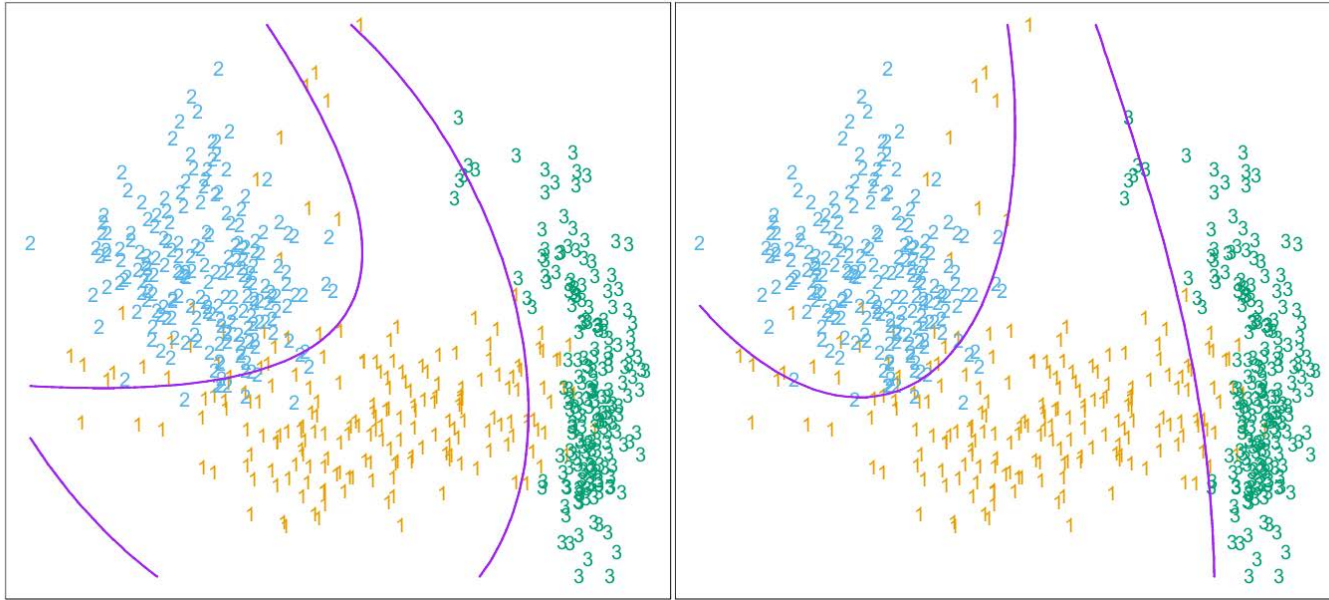
Quadratic Discriminant Analysis (QDA): The covariance matrices Σ_k are taken to be distinct (no equality simplifying assumption)

$$f_k(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)}$$

Examining the log-odds for classes k and ℓ for QDA

$$\begin{aligned} \log \frac{P(G = k | \mathbf{X} = \mathbf{x})}{P(G = \ell | \mathbf{X} = \mathbf{x})} &= \log \frac{\pi_k}{\pi_\ell} + \frac{1}{2} \log \frac{\det \Sigma_\ell}{\det \Sigma_k} - \frac{1}{2} (\mathbf{x} + \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \\ &\quad - \frac{1}{2} (\mathbf{x} + \boldsymbol{\mu}_\ell)^T \Sigma_\ell^{-1} (\mathbf{x} - \boldsymbol{\mu}_\ell) \end{aligned}$$

- There is no simplification as in LDA
- QDA produces quadratic nonlinear decision boundaries,
- The number of parameters to be estimated is higher in QDA than in LDA



Example: Three classes are utilized. The left plot shows the quadratic decision boundaries obtained using LDA by expanding the observation space to five dimensions: $X_1, X_2, X_1X_2, X_1^2, X_2^2$. The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.



Fisher Approach: Fisher derived LDA without invoking the Gaussian assumption. He posed the problem:

Find the linear combination $Z = \mathbf{a}^T \mathbf{X}$ such that the between-class variance is maximized relative to the within-class variance

Between-class variance: variance of the class means

Within-class variance: pooled variance about the class means

For \mathbf{X} , define the within-class covariance, \mathbf{W} , and between-class covariance, \mathbf{B} , matrices

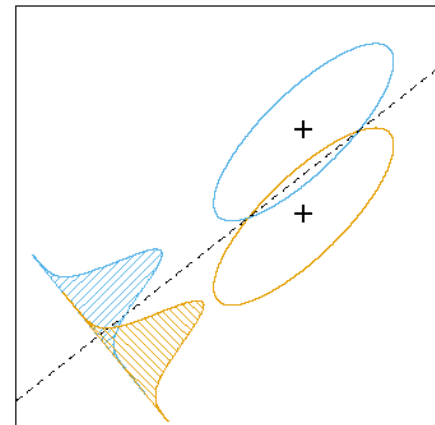
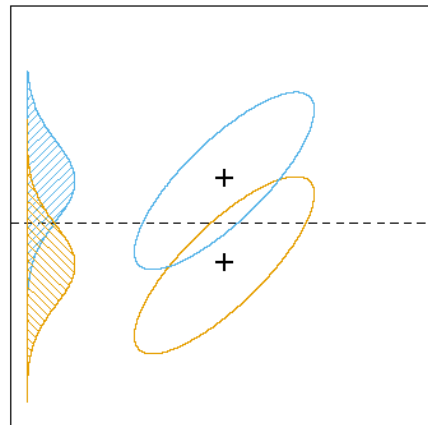
Then the between-class variance of Z is $\mathbf{a}^T \mathbf{B} \mathbf{a}$ and the within-class variance is $\mathbf{a}^T \mathbf{W} \mathbf{a}$

Fisher's formulation amounts to maximizing the *Rayleigh quotient*

$$\max_{\mathbf{a}} \frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}}$$

Or equivalently $\max_{\mathbf{a}} \mathbf{a}^T \mathbf{B} \mathbf{a}$ subject to $\mathbf{a}^T \mathbf{W} \mathbf{a} = 1$

Note: this is a generalized eigenvalue problem, with \mathbf{a} given by the largest eigenvalue of $\mathbf{W}^{-1} \mathbf{B}$



Note:

- The line through the two centroids defines the direction of greatest centroid variance
- The data projected onto that line overlaps
- The Fisher discrimination direction minimizes the overlap of projected data (Gaussian case)



Summary:

- **Classification** is the prediction of **qualitative** responses
- **Linear methods for classification** result in **linear decision boundaries**
- LDA utilizes the class posteriors $\Pr(G|\mathbf{X})$ for optimal classification and a Gaussian distribution model
 - Yields a linear boundary function, $\beta_0 + \mathbf{x}^T \boldsymbol{\beta}$
- Fisher derived the LDA result without invoking the Gaussian assumption
 - Find the linear combination $Z = \mathbf{a}^T \mathbf{X}$ such that the between-class variance is maximized relative to the within-class variance