

LDA and QDA

Eric Roemmele

Feb 15, 2018

- Let $\hat{G} : \mathcal{X} \rightarrow \mathcal{G}$ be a classifier, where \mathcal{X} is the feature space and \mathcal{G} is $\mathcal{G} = \{1, 2, \dots, K\}$ are the set of numeric class labels.
- Under zero-one loss, we can show (ESL 2.4) the optimal classifier is

$$\hat{G}(X) = k \text{ if } P(k|X = x) = \max_g P(g|X = x).$$

In other words, classify the observation to the largest posterior probability, given the features $X = x$.

- This is called the Bayes classifier, since we are looking at posterior membership probabilities.
- Most classifiers, seek to model the posterior memberships $P(G = g|X)$. (Why we can take X as “fixed”)

LDA Setup

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- Suppose the class-conditioned density $X|G = k \sim f_k(x)$, and let $\pi_k = P(G = k)$ be the prior probability of class k .
- By Bayes Theorem,

$$P(G = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{l=1}^K f_l(x)\pi_l} \quad (1)$$

- Now further suppose that $f_k(x)$ is the Gaussian density

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

- **LDA Assumption** We further suppose all classes have a common covariance matrix $\Sigma_k = \Sigma$ for $\forall k$.

Classification Boundary

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- Then, classify observation x to k if it has the largest posterior probability $P(G = k|X = x)$.
- What does our classification rule look like?
- Examine the log-ratio of two classes k and l :

$$\log \frac{P(G = k|X = x)}{P(G = l|X = x)} = \log \frac{\pi_k}{\pi_l} - \frac{1}{2}(\mu_k + \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) + x^T \Sigma^{-1}(\mu_k - \mu_l)$$

which is linear in x .

- Therefore, our decision boundaries are (linear , technically affine) p -dimensional hyperplanes.

Classification

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- Now, write

$$\begin{aligned}\log \frac{P(G = k|X = x)}{P(G = l|X = x)} > 0 &\iff x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k \\ &> x^T \Sigma^{-1} \mu_l - \frac{1}{2} \mu_l^T \Sigma^{-1} \mu_l + \log \pi_l \\ &\iff \delta_k(x) > \delta_l(x)\end{aligned}$$

- Equivalently, we can say our decision rule is $\hat{G}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$
- We call the $\delta_k(x)$ the linear discriminant functions.
- In essence, we assume we have a mixture of Gaussian distributions, and we construct the optimal linear hyperplane that separates those mixing populations.
- In practice, we need to estimate the parameters (mixing proportions, class means, and covariance matrix) by in-sample estimates.

Graph

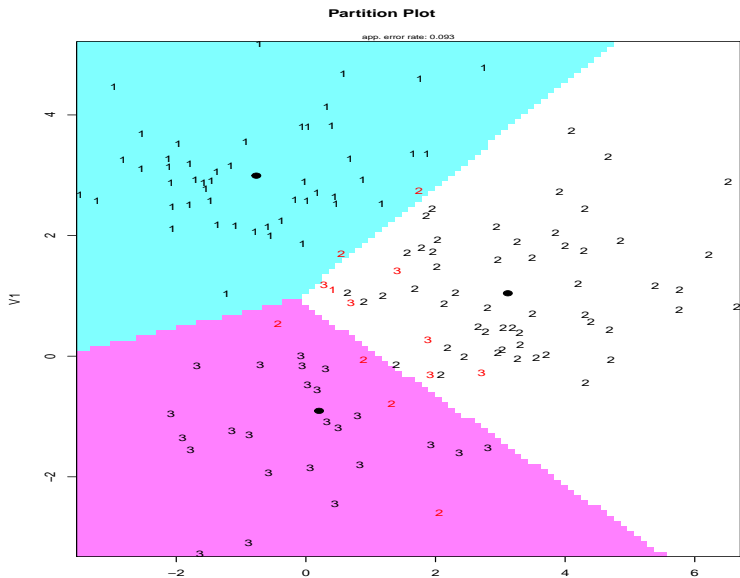
LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant



QDA

LDA and QDA

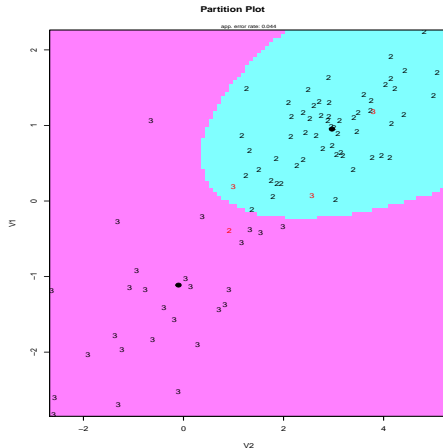
- Same process for QDA, except each group has their own covariance matrix.
- The discriminant functions are now

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$$

LDA

Performance

Fisher's
Discriminant



Why LDA?

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- LDA is surprisingly hard to beat (in terms of accuracy)
- Why does it work?
 - Our data probably aren't Gaussian
 - In the case of LDA, the covariance matrices probably aren't equal.
 - "linear discriminant analysis frequently achieves good performances in the tasks of face and object recognition, even though the assumptions of common covariance matrix among groups and normality are often violated" (Duda, et al., 2001) (Tao Li, et al., 2006)
- The reason is likely to be stability and simplicity of linear boundary.
 - More bias, but much lower variance than exotic alternatives.

Dimension Reduction in LDA

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- Another reason for LDA popularity is it contains a rank-reducing element.
- R.A. Fisher originally derived the linear classification statistic using an entirely different argument, without reference to multivariate Gaussian distributions.
- His argument was to examine linear combinations of multivariate observations x (i.e. $a^T x$) to univariate observations y (i.e. $y = a^T x$) such that the derived y from two populations were separated as much as possible.
- In other words, project the observations onto the axis where they are most separated, and then classify.

Visual

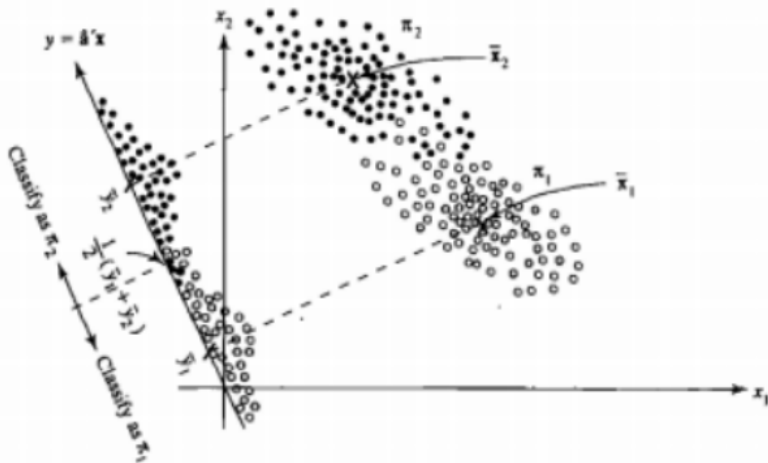
LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's Discriminant



Formulation

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- How do we find this projection?
- Define the following :

$$\text{Between Class Variance } \mathbf{B} = \sum_{i=1}^K (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T$$

$$\text{Within Class Variance } \mathbf{W} = \sum_{i=1}^K \sum_{j=1}^{n_i} (x_{ij} - \mu_i)(x_{ij} - \mu_i)^T$$

- Note that $\mathbf{S} = \mathbf{B} + \mathbf{W}$, where S is typical covariance matrix (without the (n-1) in the denominator).

Objective

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- Goal : Maximize the following function :

$$J(a) = \frac{a^T B a}{a^T W a} \quad (2)$$

- In words -*Find the linear combination $Y = a^T x$ such that the between class variance is maximized relative to the within-class variance.*
- This objective function makes sense because we want the class means well separated (i.e. high variability) relative to the within group variability.
- Thus we have the means of the groups well-separated, as well as low variability within-group. This makes classification easier, along with low dimensional representation of the data.
- The solution to (2) above is given by $a = e_1$, where e_1 is the eigenvector with the largest eigenvalue λ_1 of the matrix $W^{-1}B$. The max to (2) is λ_1 . This is not hard to prove, just use Lagrange multipliers.

Fisher's Discriminate Continued

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- We can continue with the setup in (2) and ask, what is the vector a_2 that maximizes $J(a)$ such that it is projected in the orthogonal direction to a_1 ?
- It turns out the solution is e_2 , or eigenvector of $\mathbf{W}^{-1}\mathbf{B}$ with the second largest eigenvalue.
- Usually we look for reduction of the data to two dimensions.
- Suppose we stopped in r dimensional space.

Fisher's Classification Rule

$$\hat{G}(x) = \underset{k}{\operatorname{argmin}} \sum_{j=1}^r \left[e_j^T (x - \mu_k) \right]^2$$

- To develop the above will takes a bit more work. Please refer to Johnson and Wichern Chapter 11, or any multivariate statistics book such as Rencher and Christensen.

Further LDA Related Methods to Investigate

LDA and QDA

Eric
Roemmele

LDA

Performance

Fisher's
Discriminant

- Friedman (1989) proposed a compromise between LDA and QDA where $\Sigma_k(\alpha) = \alpha \Sigma_k + (1 - \alpha) \Sigma$. Shrinking towards a common covariance matrix.
- Naive or “Idiot” Bayes - Instead of f_k the Gaussian density, replace it by a non-parametric estimate of the density assuming conditional independence between components (i.e. $f_k(x) = \prod_{i=1}^P f_{ki}(x_i)$).
- Kernel LDA - Map data to higher dimensional space, and then classify it.