



Discriminant Analysis

Daniel Khashabi¹
KHASHAB2@ILLINOIS.EDU

0.1 Introduction

Discriminant Analysis methods are among methods for linear classification or the dimensionality reduction before classification method. In LDA we assume that conditional probability densities, $f(x|Y = y)$ are normally distributed. Let's assume that for classification posterior distribution we have the followings:

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

where π_k s are class membership priors. These priors show the probability of membership in each of the classes, without having any training data at hand. Assuming multivariate Gaussian distribution,

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

The decision rule for deciding what class the data lies in is simply finding the maximum value of the class membership posterior:

$$G(x) = \arg \max_k \Pr(G = k|X = x). \quad (1)$$

In LDA we assume that all classes have the same covariance matrix $\Sigma_k = \Sigma$, $\forall k$. It can easily be shown that one can write the following expression about the log-ratio of two classes,

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

¹This is part of my notes; to see the complete list of notes check web.engr.illinois.edu/khashab2/learn.html. This work is licensed under a Creative Commons Attribution-NonCommercial 3 License.

which is the class margin, for classes l and k , when we have $\Pr(G = l|X = x) = \Pr(G = k|X = x)$. Based on the above equation we can easily see that the margin's equations are linear in terms of \mathbf{x} . That is why this method is called "Linear" Discriminant Analysis. By inspection it can be found that the decision rule in equation 1, can be equivalently written as the following equation:

$$G(x) = \arg \max_k \left\{ \mathbf{x}^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k \right\}.$$

Note that the approximate values of parameters for each class can be found using training data inside that class:

- Prior probability of data, being in one class, can be estimated using the the number of data in that class over the number of the whole data: $\hat{\pi}_k = \frac{N_k}{N}$.
- Class distribution mean can be found by averaging over the data inside the class: $\mu_k = \frac{1}{N_k} \sum_{g_i=k} \mathbf{x}_k$.
- Class covariance can be found by calculating the covariance the data inside that class, but because we have assumed in the LDA that the covariance matrix of the all classes are the same, we average the covariance matrices over all the classes, which gives us the *pooled covariance matrix*: $\hat{\Sigma} = \frac{1}{N-K} \sum_{k=1}^K \sum_{g_i=k} (\mathbf{x}_i - \hat{\mu}_k) (\mathbf{x}_i - \hat{\mu}_k)^T$

The case of *Regularized Discriminant Analysis*(RDA) is a compromise between LDA and QDA. In fact by performing an averaging on LDA and QDA's covariance matrices, we shrink the covariance of QDA towards LDA's. We can define the regularized covariance matrix as following:

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma}_k, \quad 0 \leq \alpha \leq 1, \quad (3)$$

in which the parameter α controls the shrinkage level, and RDA's inclination towards LDA or QDA, and $\hat{\Sigma}_k$ is the *pooled covariance matrix*. Now the quadratic discriminant function can be defined using the shrunken covariance matrix $\hat{\Sigma}$. The optimum value of α could be evaluated using cross-validation on test data. Another modification to the regularized version would be to shrink the covariance matrix towards a scalar matrix, which is a constant times an identity matrix:

$$\hat{\Sigma}(\gamma) = \gamma \hat{\Sigma} + (1 - \gamma) \hat{\sigma}^2 \mathbf{I}, \quad 0 \leq \gamma \leq 1, \quad (4)$$

in which γ is a tuning parameter and could be determined using cross-validation on test data. To combine two regularization parameters, simply we replace $\hat{\Sigma}$ in 3 with $\hat{\Sigma}(\gamma)$, and name the result $\hat{\Sigma}(\alpha, \gamma)$.