

Chapter 4: Linear Methods for Classification

Sections 4.1-4.3

Rebecca Steorts

Department of Statistics
University of Florida

Statistical Learning

February 1, 2010

Outline

- 1 Introduction
 - General Setup
- 2 Linear Regression
 - Naive Method
 - Iris Example
- 3 LDA
 - Assumptions
 - Derivations
 - QDA
 - Iris Example
- 4 Extensions
 - RDA
 - Reduced-Rank LDA
- 5 Conclusions

General Setup

Response categories are coded as an indicator variable. Suppose \mathcal{G} has K classes, then \mathbf{Y}_1 is a vector of 0's and 1's indicating for example whether each person is in class 1.

- The indicator response matrix is defined as $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_K)$.
- \mathbf{Y} is a matrix of 0's and 1's with each row having a single 1 indicating a person is in class k .
- The i^{th} person of interest has covariate values $\mathbf{x}_{i1}, \dots, \mathbf{x}_{ip}$ that will be represented by $\mathbf{X}_{N \times p}$.
- Our goal is to predict what class each observation is in given its covariate values.

Linear Regression of an Indicator Matrix

- Let's proceed blindly and use a naive method of linear regression.
- Fit a linear regression to each column of Y .
- The coefficient matrix is $\hat{B} = (X'X)^{-1}X'Y$.
- $\hat{Y} = X(X'X)^{-1}X'Y$
- The k^{th} column of \hat{B} contains the estimates corresponding to the linear regression coefficients that we get from regressing X_1, \dots, X_p onto Y_K .

Linear Regression of an Indicator Matrix

Look at \hat{Y} corresponding to the indicator variable for each class k . Assign each person to the class for which \hat{Y} is the largest.

More formally stated, a new observation with covariate \mathbf{x} is classified as follows:

- Compute the fitted output $\hat{\mathbf{Y}}_{new}(\mathbf{x}) = [(1, \mathbf{x})' \hat{\mathbf{B}}]'$.
- Identify the largest component of $\hat{\mathbf{Y}}_{new}(\mathbf{x})$ and classify according to

$$\hat{G}(\mathbf{x}) = \arg \max_k \hat{\mathbf{Y}}_{new}(\mathbf{x}).$$

Does this Approach Make Sense?

- The regression line estimates $E(Y_k|\mathbf{X} = \mathbf{x}) = P(G = k|\mathbf{X} = \mathbf{x})$ so the method seems somewhat sensible at first.
- Although $\sum_k \hat{Y}_k(\mathbf{x}) = 1$ for any \mathbf{x} as long as there is an intercept in the model (exercise), $\hat{Y}_k(\mathbf{x})$ can be negative or greater than 1 which is nonsensical to the initial problem statement.
- Worse problems can occur when classes are masked by others due to the rigid nature of the regression model.

Iris Data

- This data set (Fisher, Annals of Eugenics, 1936) gives the measurements of sepal and petal length and width for 150 flowers using 3 species of iris (50 flowers per species).
- The species considered are setosa, versicolor, and virginica.
- To best illustrate the methods of classification, we considered how petal width and length predict the species of a flower.

Iris Data

Sepal L	Sepal W	Petal L	Petal W	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
⋮	⋮	⋮	⋮	⋮
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.9	3.1	4.9	1.5	versicolor
⋮	⋮	⋮	⋮	⋮
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica
7.1	3.0	5.9	2.1	virginica

Illustration of Masking Effect

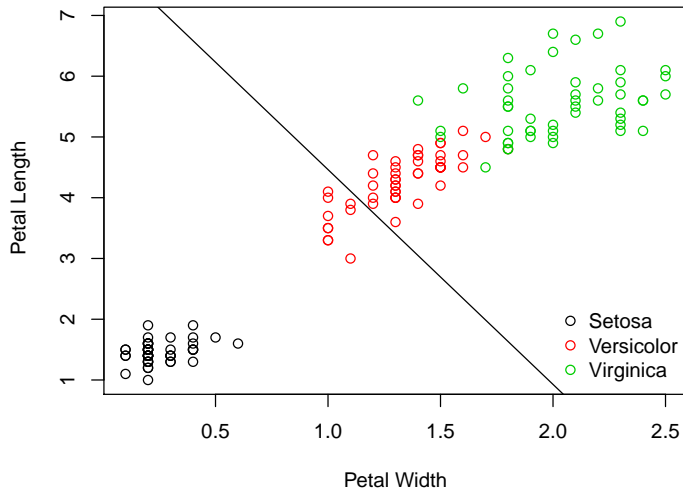


Illustration of Masking Effect

To recap the previous picture, we can see that using linear regression to predict for different classes can lead to a masking effect of one group or more. This occurs for the following reasons:

- 1 There is a plane that is high in the bottom left corner (setosa) and low in the top right corner (virginica).
- 2 There is a second plane that is high in the top right corner (virginica) but low in the bottom left corner (setosa).
- 3 The third plane is approximately flat since it tries to linearly fit a collection of points that is high in the middle (versicolor) and low on both ends.

LDA Model Assumptions

For each person, conditional on them being in class k , we assume $\mathbf{X}|G = k \sim N_p(\boldsymbol{\mu}_k, \Sigma_k)$. That is,

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}.$$

Linear Discriminant Analysis (LDA) assumes $\Sigma_k = \Sigma$ for all k .

LDA Model Assumptions

In practice the parameters of the Gaussian distribution are unknown and must be estimated by:

- $\hat{\pi}_k = N_k/N$, where N_k is the number of people of class k
- $\hat{\mu}_k = \sum_{i:g_i=k} \mathbf{x}_i / N_k$
- $\hat{\Sigma} = \sum_{k=1}^K \sum_{i:g_i=k} (\mathbf{x}_i - \hat{\mu}_k)(\mathbf{x}_i - \hat{\mu}_k)' / (N - K)$,

where $\pi_k = P(G = k)$.

Derivations

We're interested in computing

$$\begin{aligned} P(G = k | \mathbf{X} = \mathbf{x}) &= \frac{P(G = k, \mathbf{X} = \mathbf{x})}{P(\mathbf{X} = \mathbf{x})} \\ &= \frac{P(\mathbf{X} = \mathbf{x} | G = k) P(G = k)}{\sum_{k=1}^K P(\mathbf{X} = \mathbf{x}, G = k)} \\ &= \frac{f_k(\mathbf{x}) \pi_k}{\sum_{j=1}^K f_j(\mathbf{x}) \pi_j}. \end{aligned}$$

Derivations

We will compute $P(G = k | \mathbf{X} = \mathbf{x})$ for each class k .

Consider comparing $P(G = k_1 | \mathbf{X} = \mathbf{x})$ and $P(G = k_2 | \mathbf{X} = \mathbf{x})$.

Then

$$\log \left[\frac{P(G = k_1 | \mathbf{X} = \mathbf{x})}{P(G = k_2 | \mathbf{X} = \mathbf{x})} \right] = \log \left[\frac{f_{k_1}(\mathbf{x})\pi_{k_1}}{f_{k_2}(\mathbf{x})\pi_{k_2}} \right]$$

$$= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{k_1})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{k_1}) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{k_2})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{k_2}) + \log \left[\frac{\pi_{k_1}}{\pi_{k_2}} \right]$$

$$= (\boldsymbol{\mu}_{k_1} - \boldsymbol{\mu}_{k_2})' \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_{k_1}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{k_1} + \frac{1}{2} \boldsymbol{\mu}_{k_2}' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_{k_2} + \log \left[\frac{\pi_{k_1}}{\pi_{k_2}} \right]$$

Derivations

Now let's consider the boundary between predicting someone to be in class k_1 or class k_2 . To be on the the boundary, we must decide what \mathbf{x} would need to be if we think that a person is equally likely to be in class k_1 or k_2 .

This reduces to solving

$$(\mu_{k_1} - \mu_{k_2})' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{k_1}' \Sigma^{-1} \mu_{k_1} + \frac{1}{2} \mu_{k_2}' \Sigma^{-1} \mu_{k_2} + \log \left[\frac{\pi_{k_1}}{\pi_{k_2}} \right] = 0,$$

which is linear in \mathbf{x} .

- The boundary will be a line for two dimensional problems.
- The boundary will be a hyperplane for three dimensional problems.

Derivations

The linear log-odds function implies that our decision boundary between classes k_1 and k_2 will be the set where

$$P(G = k_1 | \mathbf{X} = \mathbf{x}) = P(G = k_2 | \mathbf{X} = \mathbf{x}),$$

which is linear in \mathbf{x} . In p dimensions, this is a hyperplane.

We can then say that class k_1 is more likely than class k_2 if

$$P(G = k_1 | \mathbf{X} = \mathbf{x}) > P(G = k_2 | \mathbf{X} = \mathbf{x}) \implies$$

$$\log \left[\frac{P(G = k_1 | \mathbf{X} = \mathbf{x})}{P(G = k_2 | \mathbf{X} = \mathbf{x})} \right] > 0 \implies$$

Derivations

$$(\mu_{k_1} - \mu_{k_2})' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{k_1}' \Sigma^{-1} \mu_{k_1} + \frac{1}{2} \mu_{k_2}' \Sigma^{-1} \mu_{k_2} + \log \left[\frac{\pi_{k_1}}{\pi_{k_2}} \right] > 0 \implies$$

$$\mu_{k_1}' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{k_1}' \Sigma^{-1} \mu_{k_1} + \log(\pi_{k_1}) > \mu_{k_2}' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \mu_{k_2}' \Sigma^{-1} \mu_{k_2} + \log(\pi_{k_2}).$$

Definition

The linear discriminant function $\delta_k^L(\mathbf{x})$ is defined as

$$\delta_k^L(\mathbf{x}) = \mu_k' \Sigma^{-1} \mathbf{x} - \mu_k' \Sigma^{-1} \mu_k + \log(\pi_k).$$

We can tell which class is more likely for a particular value of \mathbf{x} by comparing the classes' linear discriminant functions.

QDA

- If the Σ_k are not assumed to be equal, then convenient cancellations in our derivations earlier do not occur.
- The quadratic pieces in \mathbf{x} end up remaining leading to quadratic discriminant functions (QDA).
- QDA is similar to LDA except a covariance matrix must be estimated for each class k .

Definition

The quadratic discriminant function $\delta_k^Q(\mathbf{x})$ is defined as

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

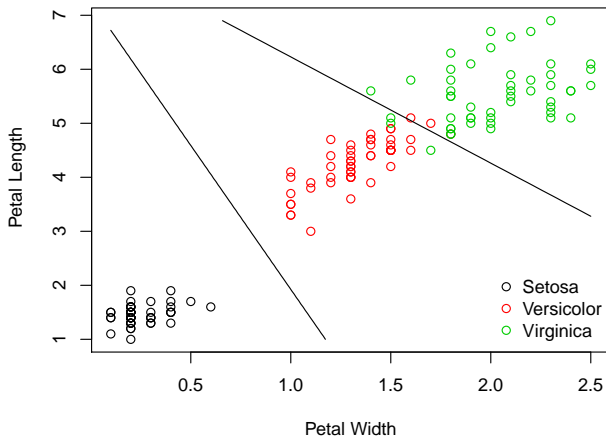
Properties of LDA and QDA

LDA and QDA seem to be widely accepted due to a bias variance trade off that leads to stability of the models.

That is, we want our model to have low variance, so we are willing to sacrifice some bias of a linear decision boundary in order for our model to be more stable.

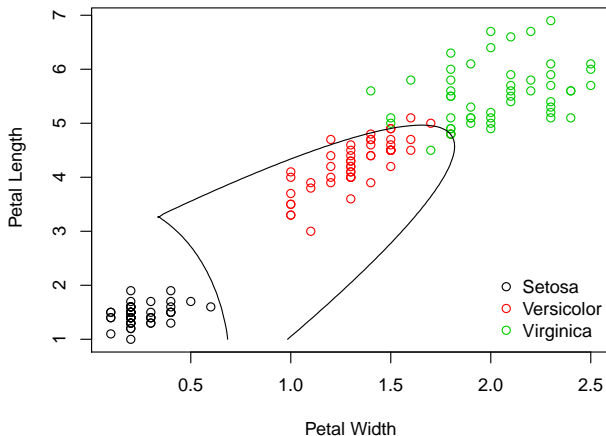
Iris Data Revisited

Returning to the iris data, we now consider predicting what classes each flower will be in using LDA. The following plot is obtained.



Iris Data Revisited

We next consider predicting what class each iris will be put into using QDA. The following plot is obtained.



Regularized Discriminant Analysis

- Friedman (1989) proposed a compromise between LDA and QDA.
- This method says that we should shrink the covariance matrices of QDA toward a common covariance matrix as done in LDA.
- Regularized covariance matrices take the form

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

- In practice, α is chosen based on performance of the model on validation data or by using cross-validation.

Geometric Interpretation of LDA

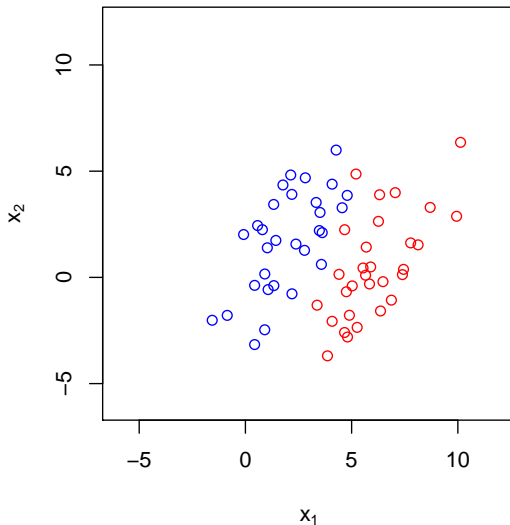
Recall K is the number of classes. Consider the centers of the K classes (centroids) which lie in a subspace of dimension $K-1$.

- We have a new observation, \mathbf{x}_{new} . For simplicity, consider $p=2$ and $K=2$.
- \mathbf{x}_1'' is defined to be the line that goes through the centroids after transforming the data to make it more spherical (ensures the line connecting the centroids is also the line that gives the greatest separation between the classes).
- \mathbf{x}_2'' is defined to be the line that is perpendicular to \mathbf{x}_1'' .
- After the data are transformed we will denote the new axes by \mathbf{x}_1' and \mathbf{x}_2' .

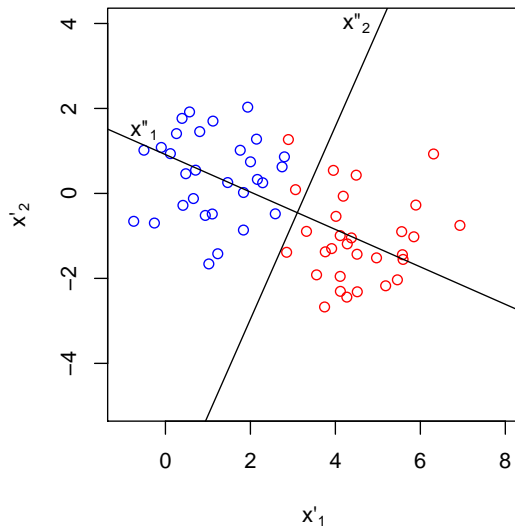
Procedure

- 1 Estimate the within class covariance matrix, $\hat{\Sigma} =: W$.
- 2 Transform the data and for each class by multiplying each data point by $W^{-1/2}$. Also, transform \mathbf{x}_{new} .
- 3 Plot \mathbf{x}' 's instead of \mathbf{x} 's to get new coordinate space.
- 4 To classify \mathbf{x}_{new} , classify this point to the class for which \mathbf{x}_{new} is closest to that class's centroid.
- 5 \mathbf{x}_1'' is the line passing through the centroids and \mathbf{x}_2'' is the line perpendicular to \mathbf{x}_1'' in the $p=2$, $K=2$ case.

Two Dimensional Example



Two Dimensional Example



Fisher's Method

Fisher proposed an alternative derivation to dimension reduction in LDA that is equivalent to the ideas previously discussed. He suggested the proper way to rotate the coordinate axes was by maximizing the variance between classes relative to the variance within the classes.

- Let $Z = \mathbf{a}'X$ and find the l.c. Z such that the between class variance is maximized wrt within class variance.
- Denote the covariance of the centroids by B .
- Denote the pooled within class covariance of the original data by W .
- BC $\text{Var}(Z) = \mathbf{a}'B\mathbf{a}$ and WC $\text{Var}(Z) = \mathbf{a}'W\mathbf{a}$.
- $B + W = T$ = total covariance matrix of X .

Fisher's Method

Fisher's problem amounts to maximizing

$$\max_{\mathbf{a}} \frac{\mathbf{a}' B \mathbf{a}}{\mathbf{a}' W \mathbf{a}} \quad (\text{exercise}).$$

The solution is \mathbf{a} = largest eigenvalue of $W^{-1}B$.

Once we find the solution to maximization problem above, denoted by \mathbf{a}_1 , we repeat the process again of maximization except this time the new maximum, \mathbf{a}_2 , must be orthogonal to \mathbf{a}_1 . This process continues and \mathbf{a}_k are called the discriminant coordinates. In terms of what we did earlier, the \mathbf{a}_k 's are equivalent to the \mathbf{x}_k'' 's.

Summary

- Using linear regression to predict for indicator variables is a naive method that can lead to entire classes being masked.
- LDA, QDA, and Regularized Discriminant Analysis procedures avoid this problem by approaching the situation with a more sound theoretical justification.
- Reduced-Rank LDA is desirable in high dimensions since it leads to a further dimension reduction in LDA.