

Chapter 12

Linear Classification

In this chapter we discuss parametric classification, in particular, linear classification, from several different points of view. We begin with a review of basic classification problems.

12.1 The Classification Problem

The problem of predicting a discrete random variable Y from another random variable X is called *classification*, also sometimes called *discrimination*, *pattern classification* or *pattern recognition*. We observe iid data $(x_1, y_1), \dots, (x_n, y_n) \sim P$ where $x_i \in \mathbb{R}^d$ and $y_i \in \{0, 1, \dots, K - 1\}$. Often, the covariates X are also called *features*. The goal is to predict Y given a new X ; here are some examples:

1. The Iris Flower study. The data are 50 samples from each of three species of Iris flowers, *Iris setosa*, *Iris virginica* and *Iris versicolor*; see Figure 12.1. The length and width of the sepal and petal are measured for each specimen, and the task is to predict the species of a new Iris flower based on these features.
2. The Coronary Risk-Factor Study (CORIS). The data consist of attributes of 462 males between the ages of 15 and 64 from three rural areas in South Africa. The outcome Y is the presence ($Y = 1$) or absence ($Y = 0$) of coronary heart disease and there are 9 covariates: systolic blood pressure, cumulative tobacco (kg), ldl (low density lipoprotein cholesterol), adiposity, famhist (family history of heart disease), typea (type-A behavior), obesity, alcohol (current alcohol consumption), and age. The goal is to predict Y from all these covariates.
3. Handwriting Digit Recognition. Here each Y is one of the ten digits from 0 to 9. There are 256 covariates X_1, \dots, X_{256} corresponding to the intensity values of the pixels in a 16×16 image; see Figure 12.2.
4. Political Blog Classification. A collection of 403 political blogs were collected during

two months before the 2004 presidential election. The goal is to predict whether a blog is *liberal* ($Y = 0$) or *conservative* ($Y = 1$) given the content of the blog.



Figure 12.1. Three different species of the Iris data. *Iris setosa* (Left), *Iris versicolor* (Middle), and *Iris virginica* (Right).

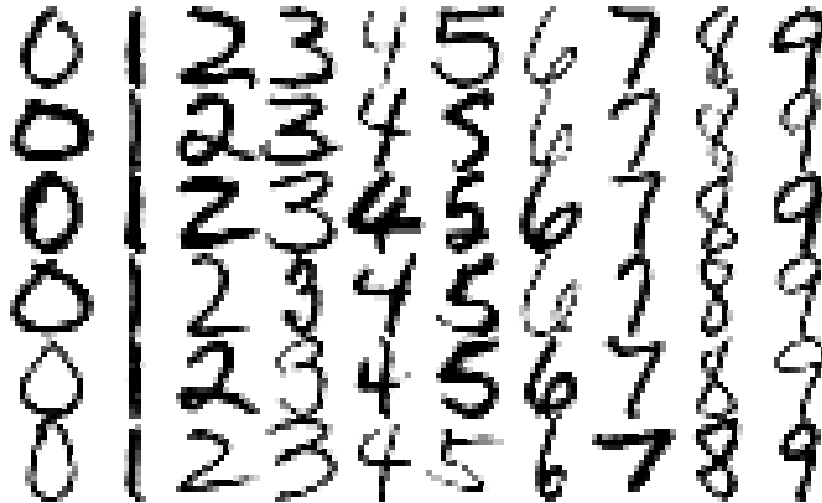


Figure 12.2. Examples from the zipcode data.

A *classification rule*, or *classifier*, is a function $h : \mathcal{X} \rightarrow \{0, \dots, K - 1\}$ where \mathcal{X} is the domain of X . When we observe a new X , we predict Y to be $h(X)$. Intuitively, the classification rule h partitions the input space \mathcal{X} into K disjoint *decision regions* whose boundaries are called *decision boundaries*. In this chapter, we mainly consider *linear classifiers* whose decision boundaries are linear functions of the covariate X . For $K = 2$, we have a *binary classification* problem. For $K > 2$, we have a *multiclass classification* problem. To simplify the discussion, we mainly discuss binary classification, and briefly explain how methods can extend to the multiclass case.

12.2 Binary Classification and Bayes Risk

A binary classifier h is a function from \mathcal{X} to $\{0, 1\}$. It is linear if there exists a function $H(x) = \beta_0 + \beta^T x$ such that $h(x) = I(H(x) > 0)$. $H(x)$ is also called a *linear discriminant function*. The decision boundary is therefore defined as the set $\{x \in \mathbb{R}^d : H(x) = 0\}$, which corresponds to a $(d - 1)$ -dimensional hyperplane within the d -dimensional input space \mathcal{X} .

The *classification risk*, or *error rate*, of h is defined as

$$R(h) = \mathbb{P}(Y \neq h(X)) \quad (12.1)$$

and the *empirical classification error* or *training error* is

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n I(h(x_i) \neq y_i). \quad (12.2)$$

We'll first establish some basic results about the general classification problem, and the relationship between classification and regression. Here is some notation that we will use throughout this chapter.

X	covariate (feature)
\mathcal{X}	domain of X , usually $\mathcal{X} \subset \mathbb{R}^d$
Y	response (pattern)
h	binary classifier, $h : \mathcal{X} \rightarrow \{0, 1\}$
H	linear discriminant function, $H(x) = \beta_0 + \beta^T x$ and $h(x) = I(H(x) > 0)$
m	regression function, $m(x) = \mathbb{E}(Y X = x) = \mathbb{P}(Y = 1 X = x)$
P_X	marginal distribution of X
p_j	$p_j(x) = p(x Y = j)$, the conditional density ^a of X given that $Y = j$
π_1	$\pi_1 = \mathbb{P}(Y = 1)$
P	joint distribution of (X, Y)

^aHere, X can be either discrete or continuous. For a discrete covariate, $p_j(x)$ is a conditional probability mass function. For a continuous covariate, $p_j(x)$ is a conditional probability density function.

First, we have the following fundamental result.

12.3 Theorem. *The rule h that minimizes $R(h)$ is*

$$h^*(x) = \begin{cases} 1 & \text{if } m(x) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad (12.4)$$

where $m(x) = \mathbb{E}(Y | X = x) = \mathbb{P}(Y = 1 | X = x)$ denotes the regression function.

The rule h^* is called the *Bayes rule*. The risk $R^* = R(h^*)$ of the Bayes rule is called the *Bayes risk*. The set $\{x \in \mathcal{X} : m(x) = 1/2\}$ is called the *Bayes decision boundary*.

Proof. We will show that $R(h) - R(h^*) \geq 0$. Note that

$$R(h) = \mathbb{P}(\{Y \neq h(X)\}) = \int \mathbb{P}(Y \neq h(X) | X = x) dP_X(x).$$

It suffices to show that

$$\mathbb{P}(Y \neq h(X) | X = x) - \mathbb{P}(Y \neq h^*(X) | X = x) \geq 0 \quad \text{for all } x \in \mathcal{X}. \quad (12.5)$$

Now,

$$\mathbb{P}(Y \neq h(X) | X = x) = 1 - \mathbb{P}(Y = h(X) | X = x) \quad (12.6)$$

$$= 1 - \left(\mathbb{P}(Y = 1, h(X) = 1 | X = x) + \mathbb{P}(Y = 0, h(X) = 0 | X = x) \right) \quad (12.7)$$

$$= 1 - \left(h(x)\mathbb{P}(Y = 1 | X = x) + (1 - h(x))\mathbb{P}(Y = 0 | X = x) \right) \quad (12.8)$$

$$= 1 - \left(h(x)m(x) + (1 - h(x))(1 - m(x)) \right). \quad (12.9)$$

Hence,

$$\begin{aligned} & \mathbb{P}(Y \neq h(X) | X = x) - \mathbb{P}(Y \neq h^*(X) | X = x) \\ &= \left(h^*(x)m(x) + (1 - h^*(x))(1 - m(x)) \right) - \left(h(x)m(x) + (1 - h(x))(1 - m(x)) \right) \\ &= (2m(x) - 1)(h^*(x) - h(x)) = 2 \left(m(x) - \frac{1}{2} \right) (h^*(x) - h(x)). \end{aligned} \quad (12.10)$$

When $m(x) \geq 1/2$ and $h^*(x) = 1$, (12.10) is non-negative. When $m(x) < 1/2$ and $h^*(x) = 0$, (12.10) is again non-negative. This proves (12.5). \square

We can rewrite h^* in a different way. From Bayes' theorem

$$\begin{aligned} m(x) &= \mathbb{P}(Y = 1 | X = x) = \frac{p(x | Y = 1)\mathbb{P}(Y = 1)}{p(x | Y = 1)\mathbb{P}(Y = 1) + p(x | Y = 0)\mathbb{P}(Y = 0)} \\ &= \frac{\pi_1 p_1(x)}{\pi_1 p_1(x) + (1 - \pi_1) p_0(x)}. \end{aligned} \quad (12.11)$$

where $\pi_1 = \mathbb{P}(Y = 1)$. From the above equality, we have that

$$m(x) > \frac{1}{2} \quad \text{is equivalent to} \quad \frac{p_1(x)}{p_0(x)} > \frac{1 - \pi_1}{\pi_1}. \quad (12.12)$$

Thus the Bayes rule can be rewritten as

$$h^*(x) = \begin{cases} 1 & \text{if } \frac{p_1(x)}{p_0(x)} > \frac{1 - \pi_1}{\pi_1} \\ 0 & \text{otherwise.} \end{cases} \quad (12.13)$$

If \mathcal{H} is a set of classifiers then the classifier $h_o \in \mathcal{H}$ that minimizes $R(h)$ is the *oracle classifier*. Formally,

$$R(h_o) = \inf_{h \in \mathcal{H}} R(h)$$

and $R_o = R(h_o)$ is called the *oracle risk* of \mathcal{H} . In general, if h is any classifier and R^* is the Bayes risk then,

$$R(h) - R^* = \underbrace{R(h) - R(h_o)}_{\text{distance from oracle}} + \underbrace{R(h_o) - R^*}_{\text{distance of oracle from Bayes error}}. \quad (12.14)$$

The first term is analogous to the variance, and the second is analogous to the squared bias in linear regression.

12.3 Classification is Easier than Regression

For a binary classifier problem, given a covariate X we only need to predict its class label $Y = 0$ or $Y = 1$. This is in contrast to a regression problem where we need to predict a real-valued response $Y \in \mathbb{R}$. Intuitively, classification is a much easier task than regression. To rigorously formalize this, let $m^*(x) = \mathbb{E}(Y | X = x)$ be the true regression function and let $h^*(x)$ be the corresponding Bayes rule. Let $\hat{m}(x)$ be an estimate of $m^*(x)$ and define the *plug-in classification rule*:

$$\hat{h}(x) = \begin{cases} 1 & \text{if } \hat{m}(x) > \frac{1}{2} \\ 0 & \text{otherwise.} \end{cases} \quad (12.15)$$

We have the following theorem.

12.16 Theorem. *The risk of the plug-in classifier rule in (12.15) satisfies*

$$R(\hat{h}) - R^* \leq 2 \sqrt{\int (\hat{m}(x) - m^*(x))^2 dP_X(x)}.$$

Proof. In the proof of Theorem 12.3 we showed that

$$\begin{aligned} \mathbb{P}(Y \neq \hat{h}(X) | X = x) - \mathbb{P}(Y \neq h^*(X) | X = x) &= (2\hat{m}(x) - 1)(h^*(x) - \hat{h}(x)) \\ &= |2\hat{m}(x) - 1| I(h^*(x) \neq \hat{h}(x)) = 2|\hat{m}(x) - 1/2| I(h^*(x) \neq \hat{h}(x)). \end{aligned}$$

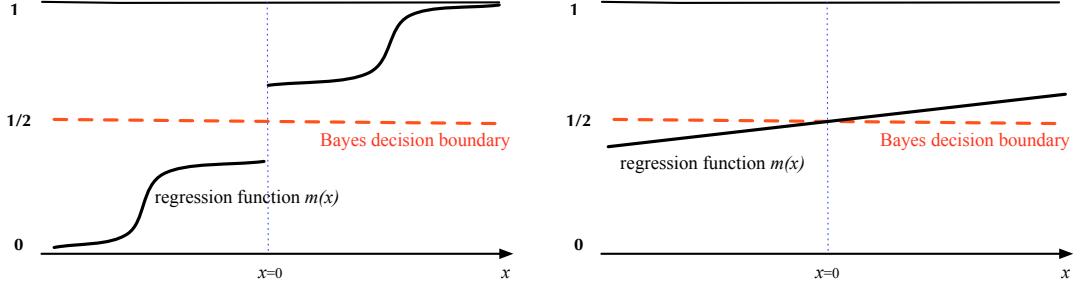


Figure 12.3. The Bayes rule is $h^*(x) = I(x > 0)$ in both plots, which show the regression function $m(x) = \mathbb{E}(Y | x)$ for two problems. The left plot shows an easy problem; there is little ambiguity around the decision boundary. The right plot shows a hard problem; it is hard to know from the data if you are to the left or right of the decision boundary.

Now, when $h^*(x) \neq \hat{h}(x)$, there are two possible cases: (i) $\hat{h}(x) = 1$ and $h^*(x) = 0$; (ii) $\hat{h}(x) = 0$ and $h^*(x) = 1$. In both cases, we have that $|\hat{m}(x) - m^*(x)| \geq |\hat{m}(x) - 1/2|$. Therefore,

$$\begin{aligned} \mathbb{P}(\hat{h}(X) \neq Y) - \mathbb{P}(h^*(X) \neq Y) &= 2 \int |\hat{m}(x) - 1/2| I(h^*(x) \neq \hat{h}(x)) dP_X(x) \\ &\leq 2 \int |\hat{m}(x) - m^*(x)| I(h^*(x) \neq \hat{h}(x)) dP_X(x) \\ &\leq 2 \int |\hat{m}(x) - m^*(x)| dP_X(x) \end{aligned} \quad (12.17)$$

$$\leq 2 \sqrt{\int (\hat{m}(x) - m^*(x))^2 dP_X(x)}. \quad (12.18)$$

The last inequality follows from the fact that $\mathbb{E}|Z| \leq \sqrt{\mathbb{E}Z^2}$ for any Z . \square

This theorem implies that if the regression estimate $\hat{m}(x)$ is close to $m^*(x)$ then the plug-in classification risk will be close to the Bayes risk. The converse is *not* necessarily true. It is possible for \hat{m} to be far from $m^*(x)$ and still lead to a good classifier. As long as $\hat{m}(x)$ and $m^*(x)$ are on the same side of $1/2$ they yield the same classifier.

12.19 Example. Figure 12.3 shows two one-dimensional regression functions. In both cases, the Bayes rule is $h^*(x) = I(x > 0)$ and the decision boundary is $\mathcal{D} = \{x = 0\}$. The left plot illustrates an easy problem; there is little ambiguity around the decision boundary. Even a poor estimate of $m(x)$ will recover the correct decision boundary. The right plot illustrates a hard problem; it is hard to know from the data if you are to the left or right of the decision boundary. \square

12.4 Parametric Models

The term “parametric classification” can mean several things. For example:

- The conditional densities $p_0(x) = p(x | Y = 0)$ and $p_1(x) = p(x | Y = 1)$ are assumed to belong to a parametric family $\mathcal{M} = \{p_\theta(x) : \theta \in \Theta\}$.
- The conditional probability $\mathbb{P}(Y = 1 | X = x)$ for Y given $X = x$ is assumed to belong to a parametric family $\{\mathbb{P}(Y = 1 | X = x; \theta) : \theta \in \Theta\}$.
- The set of classifiers \mathcal{H} is assumed to be parametric:

$$\mathcal{H} = \{h(x; \theta) : \mathcal{X} \rightarrow \{0, 1\} : \theta \in \Theta\}.$$

The first two cases are often called *generative models* and *discriminative models*. In more detail, the joint density of a single observation (x_i, y_i) is $p(x_i, y_i) = p(x_i | y_i)p(y_i) = p(y_i | x_i)p(x)$. In the generative case we often estimate the joint distribution by maximizing the *joint likelihood*:

$$\prod_{i=1}^n p(x_i, y_i) = \underbrace{\prod_{i=1}^n p(x_i | y_i)}_{\text{parametric model}} \underbrace{\prod_{i=1}^n p(y_i)}_{\text{Bernoulli}}. \quad (12.20)$$

In the discriminative case we maximize the conditional likelihood $\prod_{i=1}^n p(y_i | x_i)$ and ignore the second term $p(x_i)$:

$$\prod_{i=1}^n p(x_i, y_i) = \underbrace{\prod_{i=1}^n p(y_i | x_i)}_{\text{parametric model}} \underbrace{\prod_{i=1}^n p(x_i)}_{\text{ignored}}. \quad (12.21)$$

Since classification only requires knowing $p(y_i | x_i)$, we don’t really need to estimate the whole joint distribution. Discriminative models leave the marginal distribution $p(x)$ unspecified, and so are “more nonparametric” by relying on fewer distributional assumptions.

In the generative case, we parameterize the conditional densities $p_0(x)$ and $p_1(x)$ as $p_{\theta_0,0}(x)$ and $p_{\theta_1,1}(x)$, where $p_{\theta_0,0}(x) = p_{\theta_0}(x | Y = 0)$ and $p_{\theta_1,1}(x) = p_{\theta_1}(x | Y = 1)$. In this case, the regression function $m(x) = \mathbb{P}(Y = 1 | X = x)$ is also a function of $\theta = (\theta_1, \theta_2)^T$, we denote it by $m_\theta(x)$. Thus,

$$m_\theta(x) \equiv \mathbb{P}(Y = 1 | X = x) = \frac{\pi_1 p_{\theta_1,1}(x)}{(1 - \pi_1) p_{\theta_0,0}(x) + \pi_1 p_{\theta_1,1}(x)}.$$

Given an estimator $(\hat{\theta}_n, \hat{\pi}_1)$, we define the plug-in estimator

$$\hat{h}(x) = I(m_{\hat{\theta}_n}(x) > 1/2).$$

Often, $\hat{\theta}_n$ is the maximum likelihood estimator.

12.5 Gaussian Discriminant Analysis

Suppose that $p_0(x) = p(x | Y = 0)$ and $p_1(x) = p(x | Y = 1)$ are both multivariate Gaussians:

$$p_k(x) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right\}, \quad k = 0, 1.$$

where Σ_1 and Σ_2 are both $d \times d$ covariance matrices. Thus, $X | Y = 0 \sim N(\mu_0, \Sigma_0)$ and $X | Y = 1 \sim N(\mu_1, \Sigma_1)$.

Given a square matrix A , we define $|A|$ to be the determinant of A . For a binary classification problem with Gaussian distributions, we have the following theorem.

12.22 Theorem. *If $X | Y = 0 \sim N(\mu_0, \Sigma_0)$ and $X | Y = 1 \sim N(\mu_1, \Sigma_1)$, then the Bayes rule is*

$$h^*(x) = \begin{cases} 1 & \text{if } r_1^2 < r_0^2 + 2 \log \left(\frac{\pi_1}{1-\pi_1} \right) + \log \left(\frac{|\Sigma_0|}{|\Sigma_1|} \right) \\ 0 & \text{otherwise} \end{cases} \quad (12.23)$$

where $r_i^2 = (x - \mu_i)^T \Sigma_i^{-1} (x - \mu_i)$ for $i = 1, 2$ is the Mahalanobis distance.

Proof. By definition, the Bayes rule is $h^*(x) = I(\pi_1 p_1(x) > (1 - \pi_1) p_0(x))$. Plug-in the specific forms of p_0 and p_1 and take the logarithms we get $h^*(x) = 1$ if and only if

$$\begin{aligned} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) - 2 \log \pi_1 + \log(|\Sigma_1|) \\ < (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) - 2 \log(1 - \pi_1) + \log(|\Sigma_0|). \end{aligned} \quad (12.24)$$

The theorem immediately follows from some simple algebra. \square

Let $\pi_0 = 1 - \pi_1$. An equivalent way of expressing the Bayes rule is

$$h^*(x) = \operatorname{argmax}_{k \in \{0,1\}} \delta_k(x) \quad (12.25)$$

where

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k \quad (12.26)$$

is called the *Gaussian discriminant function*. The decision boundary of the above classifier can be characterized by the set $\{x \in \mathcal{X} : \delta_1(x) = \delta_0(x)\}$, which is quadratic. So, this procedure is called *quadratic discriminant analysis* (QDA).

In practice, we use sample quantities of $\pi_0, \pi_1, \mu_1, \mu_2, \Sigma_0, \Sigma_1$ in place of their popula-

tion values, namely (See Exercise 4):

$$\hat{\pi}_0 = \frac{1}{n} \sum_{i=1}^n (1 - y_i), \quad \hat{\pi}_1 = \frac{1}{n} \sum_{i=1}^n y_i, \quad (12.27)$$

$$\hat{\mu}_0 = \frac{1}{n_0} \sum_{i: y_i=0} x_i, \quad \hat{\mu}_1 = \frac{1}{n_1} \sum_{i: y_i=1} x_i, \quad (12.28)$$

$$\hat{\Sigma}_0 = \frac{1}{n_0 - 1} \sum_{i: y_i=0} (x_i - \hat{\mu}_0)(x_i - \hat{\mu}_0)^T, \quad (12.29)$$

$$\hat{\Sigma}_1 = \frac{1}{n_1 - 1} \sum_{i: y_i=1} (x_i - \hat{\mu}_1)(x_i - \hat{\mu}_1)^T, \quad (12.30)$$

where $n_0 = \sum_i (1 - y_i)$ and $n_1 = \sum_i y_i$. (Note: we could also estimate Σ_0 and Σ_1 using their maximum likelihood estimates, which replace $n_0 - 1$ and $n_1 - 1$ with n_0 and n_1 .)

A simplification occurs if we assume that $\Sigma_0 = \Sigma_1 = \Sigma$. In this case, the Bayes rule is

$$h^*(x) = \operatorname{argmax}_k \delta_k(x) \quad (12.31)$$

where now

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k. \quad (12.32)$$

The parameters are estimated as before, except that we use a pooled estimate of Σ :

$$\hat{\Sigma} = \frac{(n_0 - 1) \hat{\Sigma}_0 + (n_1 - 1) \hat{\Sigma}_1}{n_0 + n_1 - 2}. \quad (12.33)$$

The classification rule is

$$h^*(x) = \begin{cases} 1 & \text{if } \delta_1(x) > \delta_0(x) \\ 0 & \text{otherwise.} \end{cases} \quad (12.34)$$

The decision boundary $\{x \in \mathcal{X} : \delta_0(x) = \delta_1(x)\}$ is linear, so this method is called *linear discrimination analysis* (LDA).

When the dimension d is large, fully specifying the QDA decision boundary requires $d + d(d - 1)$ parameters, and fully specifying the LDA decision boundary requires $d + d(d - 1)/2$ parameters. Such a large number of free parameters might induce a large variance. To further regularize the model, two popular methods are *diagonal quadratic discriminant analysis* (DQDA) and *diagonal linear discriminant analysis* (DLDA). The only difference between DQDA and DLDA with QDA and LDA is that after calculating $\hat{\Sigma}_1$ and $\hat{\Sigma}_0$ as in (12.30), we set all the off-diagonal elements to be zero. This is also called the “independence rule” in Bickel and Levina (2004).

12.6 Multiclass Gaussian Discriminant Analysis

We now generalize to the case where Y takes on more than two values. That is, $Y \in \{0, \dots, K-1\}$ for $K > 2$. First, we characterize the Bayes classifier under this multiclass setting.

12.35 Theorem. *Let $R(h) = \mathbb{P}(h(X) \neq Y)$ be the classification error of a classification rule $h(x)$. The Bayes rule $h^*(X)$ minimizing $R(h)$ can be written as*

$$h^*(x) = \operatorname{argmax}_k \mathbb{P}(Y = k | X = x) \quad (12.36)$$

Proof. We have

$$R(h) = 1 - \mathbb{P}(h(X) = Y) \quad (12.37)$$

$$= 1 - \sum_{k=0}^{K-1} \mathbb{P}(h(X) = k, Y = k) \quad (12.38)$$

$$= 1 - \sum_{k=0}^{K-1} \mathbb{E} \left[I(h(X) = k) \mathbb{P}(Y = k | X) \right] \quad (12.39)$$

It's clear that $h^*(X) = \operatorname{argmax}_k \mathbb{P}(Y = k | X)$ achieves the minimized classification error $1 - \mathbb{E}[\max_k \mathbb{P}(Y = k | X)]$. \square

Let $\pi_k = \mathbb{P}(Y = k)$. The next theorem extends QDA and LDA to the multiclass setting.

12.40 Theorem. *Suppose that $Y \in \{0, \dots, K-1\}$ with $K \geq 2$. If $p_k(x) = p(x | Y = k)$ is Gaussian: $X | Y = k \sim N(\mu_k, \Sigma_k)$, the Bayes rule for the multiclass QDA can be written as*

$$h^*(x) = \operatorname{argmax}_k \delta_k(x)$$

where

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k. \quad (12.41)$$

If all Gaussians have an equal variance Σ , then

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k. \quad (12.42)$$

Proof. See Exercise 6 \square

Let $n_k = \sum_i I(y_i = k)$ for $k = 0, \dots, K-1$. The estimated sample quantities of π_k ,

μ_k , Σ_k , and Σ are:

$$\hat{\pi}_k = \frac{1}{n} \sum_{i=1}^n I(y_i = k), \quad \hat{\mu}_k = \frac{1}{n_k} \sum_{i: y_i=k} x_i, \quad (12.43)$$

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i: y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T, \quad (12.44)$$

$$\hat{\Sigma} = \frac{\sum_{k=0}^{K-1} (n_k - 1) \hat{\Sigma}_k}{n - K}. \quad (12.45)$$

12.46 Example. Let us return to the Iris data example. Recall that there are 150 observations made on three classes of the iris flower: *Iris setosa*, *Iris versicolor*, and *Iris virginica*. There are four features: sepal length, sepal width, petal length, and petal width. In Figure 12.4 we visualize the datasets. Within each class, we plot the densities for each feature. It's easy to see that the distributions of petal length and petal width are quite different across different classes, which suggests that they are very informative features.

Figures 12.5 and 12.6 provide multiple figure arrays illustrating the classification of observations based on LDA and QDA for every combination of two features. The classification boundaries and error are obtained by simply restricting the data to these a given pair of features before fitting the model. We see that the decision boundaries for LDA are linear, while the decision boundaries for QDA are highly nonlinear. The training errors for LDA and QDA on this data are both 0.02. From these figures, we see that it is very easy to discriminate the observations of class *Iris setosa* from those of the other two classes. \square

12.7 Binary Class Logistic Regression

One approach to binary classification is to estimate the regression function $m(x) = \mathbb{E}(Y|X = x) = \mathbb{P}(Y = 1|X = x)$ and, once we have an estimate $\hat{m}(x)$, use the classification rule

$$\hat{h}(x) = \begin{cases} 1 & \text{if } \hat{m}(x) > \frac{1}{2} \\ 0 & \text{otherwise.} \end{cases} \quad (12.47)$$

For binary classification problems, one possible choice is the linear regression model

$$Y = m(X) + \epsilon = \beta_0 + \sum_{j=1}^d \beta_j X_j + \epsilon. \quad (12.48)$$

Later, we will show that if Y is coded as $\{-1, +1\}$, the linear regression model is equivalent to Fisher linear discriminant analysis. Nonetheless, the linear regression model does not explicitly constrain Y to take on binary values. A more natural alternative is to use *logistic regression*, which is the most common binary classification method.

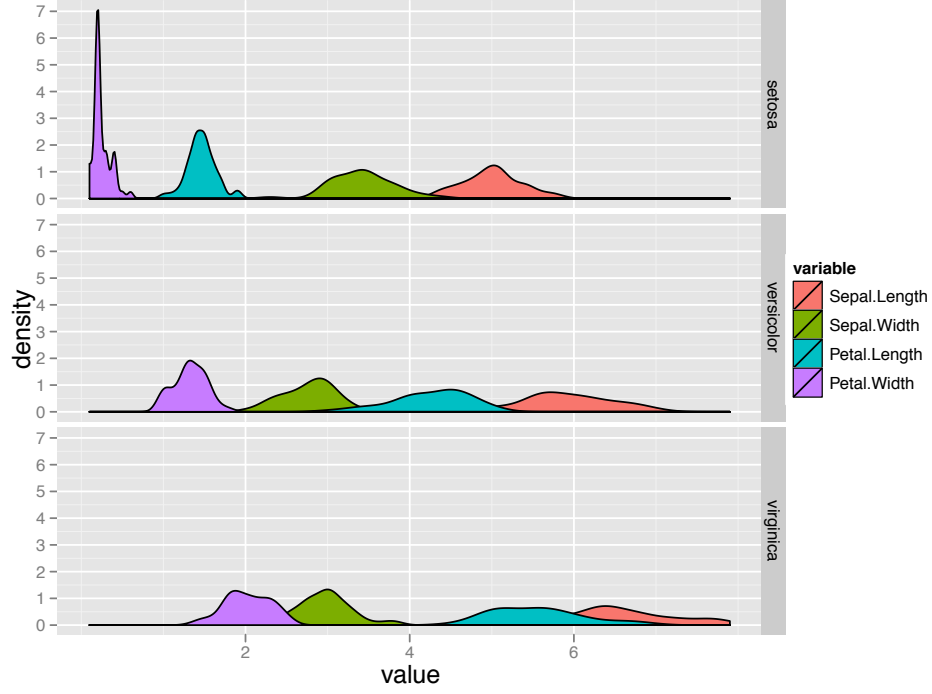


Figure 12.4. *The Iris data: The estimated densities for different features are plotted within each class. It's easy to see that the distributions of petal length and petal width are quite different across different classes, which suggests that they are very informative features.*

Before we describe the logistic regression model, let's recall some basic facts about binary random variables. If Y takes values 0 and 1, we say that Y has a Bernoulli distribution with parameter $\pi_1 = \mathbb{P}(Y = 1)$. The probability mass function for Y is $p(y; \pi_1) = \pi_1^y (1 - \pi_1)^{1-y}$ for $y = 0, 1$. The likelihood function for π_1 based on iid data y_1, \dots, y_n is

$$\mathcal{L}(\pi_1) = \prod_{i=1}^n p(y_i; \pi_1) = \prod_{i=1}^n \pi_1^{y_i} (1 - \pi_1)^{1-y_i}. \quad (12.49)$$

In the logistic regression model, we assume that

$$m(x) = \mathbb{P}(Y = 1 \mid X = x) = \frac{\exp(\beta_0 + x^T \beta)}{1 + \exp(\beta_0 + x^T \beta)} \equiv \pi_1(x, \beta_0, \beta). \quad (12.50)$$

In other words, given $X = x$, Y is Bernoulli with mean $\pi_1(x, \beta_0, \beta)$. We can write the model as

$$\text{logit}(\mathbb{P}(Y = 1 \mid X = x)) = \beta_0 + x^T \beta \quad (12.51)$$

where $\text{logit}(a) = \log(a/(1 - a))$. The name “logistic regression” comes from the fact that $\exp(x)/(1 + \exp(x))$ is called the logistic function. A plot of the logistic function for a one-dimensional covariate is shown in Figure 12.7.

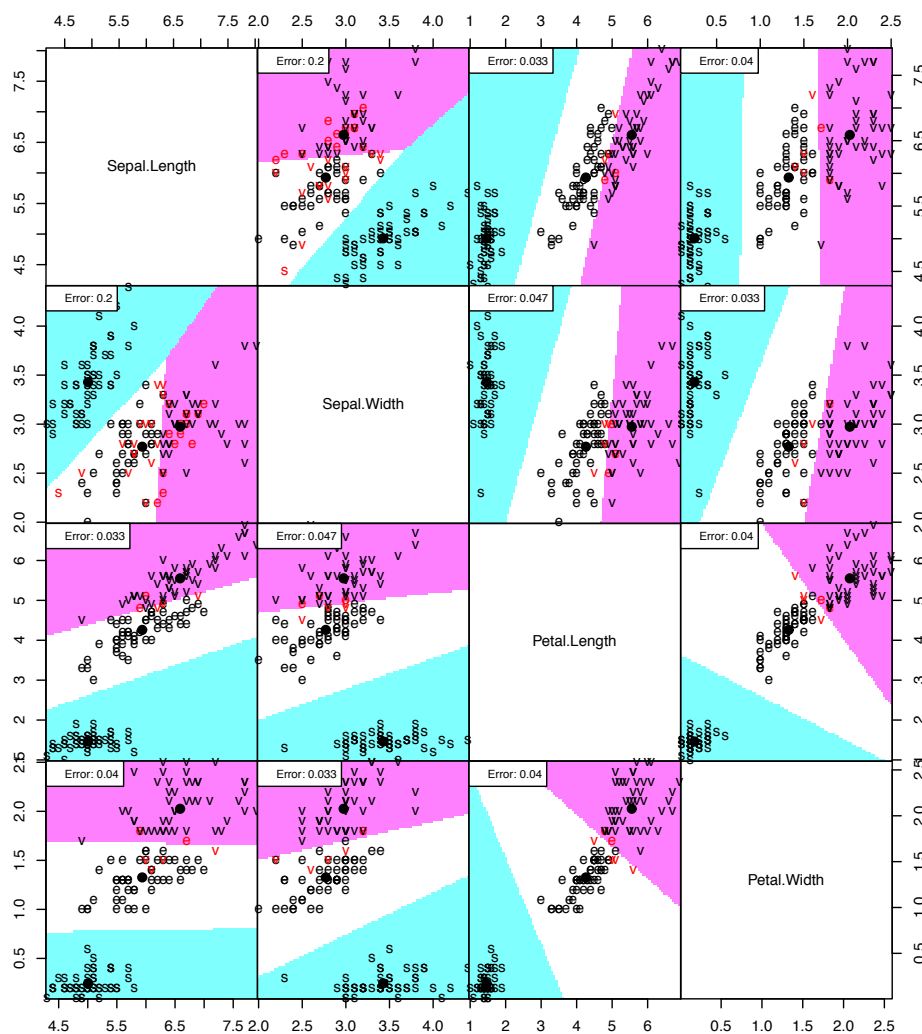


Figure 12.5. Classifying the Iris data using LDA. The multiple figure array illustrates the classification of observations based on LDA for every combination of two features. The classification boundaries and error are obtained by simply restricting the data to a given pair of features before fitting the model. In these plots, “s” represents the class label Iris setosa, “e” represents the class label Iris versicolor, and “v” represents the class label Iris virginica. The red letters illustrate the misclassified observations.

12.52 Lemma. Both linear regression and logistic regression models have linear decision boundaries.

Proof. The linear decision boundary for linear regression is straightforward. The same result for logistic regression follows from the monotonicity of the logistic function. \square

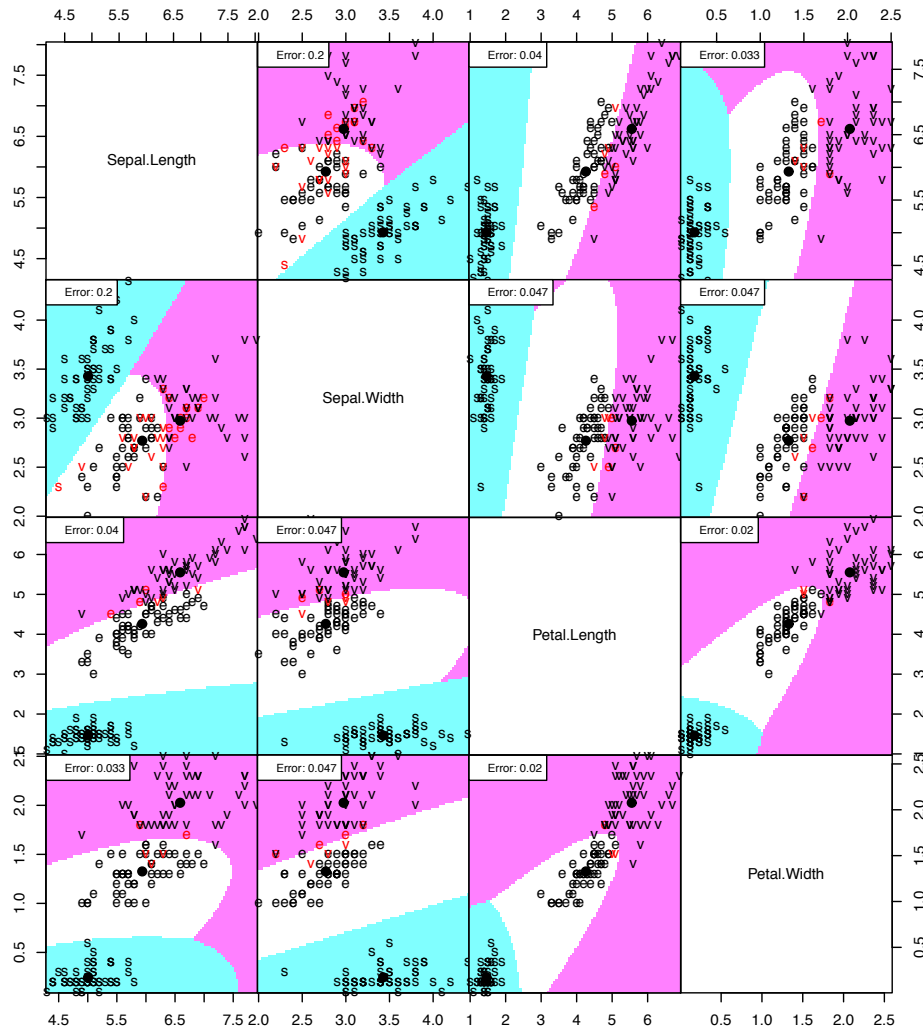


Figure 12.6. Classifying the Iris data using QDA. The multiple figure array illustrates the classification of observations based on QDA for every combination of two features. The classification boundaries are displayed and the classification error by simply casting the data onto these two features are calculated. In these plots, “s” represents the class label *Iris setosa*, “e” represents the class label *Iris versicolor*, and “v” represents the class label *Iris virginica*. The red letters illustrate the misclassified observations.

The parameters β_0 and $\beta = (\beta_1, \dots, \beta_d)^T$ can be estimated by maximum conditional likelihood. The conditional likelihood function for β is

$$\mathcal{L}(\beta_0, \beta) = \prod_{i=1}^n \pi_1(x_i, \beta_0, \beta)^{y_i} (1 - \pi_1(x_i, \beta_0, \beta))^{1-y_i}.$$

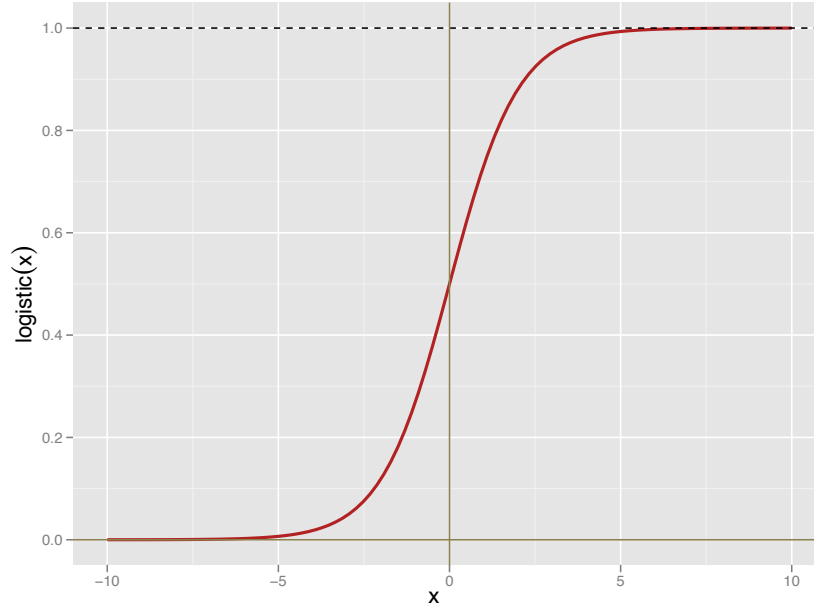


Figure 12.7. *The logistic function*

Thus the conditional log-likelihood is

$$\ell(\beta_0, \beta) = \sum_{i=1}^n \left\{ y_i \log \pi_1(x_i, \beta_0, \beta) - (1 - y_i) \log(1 - \pi_1(x_i, \beta_0, \beta)) \right\} \quad (12.53)$$

$$= \sum_{i=1}^n \left\{ y_i (\beta_0 + x_i^T \beta) - \log(1 + \exp(\beta_0 + x_i^T \beta)) \right\}. \quad (12.54)$$

The maximum conditional likelihood estimators $\hat{\beta}_0$ and $\hat{\beta}$ cannot be found in closed form. However, the loglikelihood function is concave and can be efficiently solve by the Newton's method in an iterative manner.

12.8 Fitting Binary Class Logistic Regression

For notational simplicity, we redefine (local to this section) the d -dimensional covariate x_i and parameter vector β as the following $(d + 1)$ -dimensional vectors:

$$x_i \leftarrow (1, x_i^T)^T \text{ and } \beta \leftarrow (\beta_0, \beta^T)^T. \quad (12.55)$$

Thus, we simplify $\pi_1(x, \beta_0, \beta)$ as $\pi_1(x, \beta)$ and $\ell(\beta_0, \beta)$ as $\ell(\beta)$.

To maximize $\ell(\beta)$, the $(k + 1)$ th Newton step in the algorithm replaces the k th iterate

$\hat{\beta}^{(k)}$ by

$$\hat{\beta}^{(k+1)} \leftarrow \hat{\beta}^{(k)} - \left(\frac{\partial^2 \ell(\hat{\beta}^{(k)})}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\hat{\beta}^{(k)})}{\partial \beta}. \quad (12.56)$$

The gradient $\frac{\partial \ell(\hat{\beta}^{(k)})}{\partial \beta}$ and Hessian $\frac{\partial^2 \ell(\hat{\beta}^{(k)})}{\partial \beta \partial \beta^T}$ are both evaluated at $\hat{\beta}^{(k)}$ and can be written as (see Exercise 7)

$$\frac{\partial \ell(\hat{\beta}^{(k)})}{\partial \beta} = \sum_{i=1}^n (\pi_1(x_i, \hat{\beta}^{(k)}) - y_i) x_i \text{ and } \frac{\partial^2 \ell(\hat{\beta}^{(k)})}{\partial \beta \partial \beta^T} = -\mathbf{X}^T \mathbf{W} \mathbf{X} \quad (12.57)$$

where $\mathbf{W} = \text{diag}(w_{11}^{(k)}, w_{22}^{(k)}, \dots, w_{dd}^{(k)})$ is a diagonal matrix with

$$w_{ii}^{(k)} = \pi_1(x_i, \hat{\beta}^{(k)}) (1 - \pi_1(x_i, \hat{\beta}^{(k)})). \quad (12.58)$$

Let $\pi_1^{(k)} = (\pi_1(x_1, \hat{\beta}^{(k)}), \dots, \pi_1(x_n, \hat{\beta}^{(k)}))^T$, (12.56) can be written as

$$\hat{\beta}^{(k+1)} = \hat{\beta}^{(k)} + (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \pi_1^{(k)}) \quad (12.59)$$

$$= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} (\mathbf{X} \hat{\beta}^{(k)} + \mathbf{W}^{-1} (\mathbf{y} - \pi_1^{(k)})) \quad (12.60)$$

$$= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{z}^{(k)} \quad (12.61)$$

where $\mathbf{z}^{(k)} \equiv (z_1^{(k)}, \dots, z_n^{(k)})^T = \mathbf{X}^T \hat{\beta}^{(k)} + \mathbf{W}^{-1} (\mathbf{y} - \pi_1^{(k)})$ with

$$z_i^{(k)} = \log \left(\frac{\pi_1(x_i, \hat{\beta}^{(k)})}{1 - \pi_1(x_i, \hat{\beta}^{(k)})} \right) + \frac{y_i - \pi_1(x_i, \hat{\beta}^{(k)})}{\pi_1(x_i, \hat{\beta}^{(k)}) (1 - \pi_1(x_i, \hat{\beta}^{(k)}))}. \quad (12.62)$$

Given the current estimate $\hat{\beta}^{(k)}$, the above Newton iteration forms a quadratic approximation to the negative log-likelihood using Taylor expansion at $\hat{\beta}^{(k)}$:

$$-\ell(\beta) = \underbrace{\frac{1}{2} (\mathbf{z} - \mathbf{X}\beta)^T \mathbf{W} (\mathbf{z} - \mathbf{X}\beta)}_{\ell_Q(\beta)} + \text{constant}. \quad (12.63)$$

Essentially, the update equation (12.61) corresponds to solving a quadratic optimization

$$\hat{\beta}^{(k+1)} = \arg \min_{\beta} \ell_Q(\beta). \quad (12.64)$$

We then get an iterative algorithm called *iteratively reweighted least squares*.

Iteratively Reweighted Least Squares Algorithm

Choose starting values $\hat{\beta}^{(0)} = (\hat{\beta}_0^{(0)}, \hat{\beta}_1^{(0)}, \dots, \hat{\beta}_d^{(0)})^T$ and compute $\pi_1(x_i, \hat{\beta}^{(0)})$ using Equation (12.50), for $i = 1, \dots, n$ with β_j replaced by its initial value $\hat{\beta}_j^{(0)}$.

For $k = 1, 2, \dots$, iterate the following steps until convergence.

1. Calculate $z_i^{(k)}$ according to (12.62) for $i = 1, \dots, n$.
2. Calculate $\hat{\beta}^{(k+1)}$ according to (12.61). This corresponds to doing a weighted linear regression of \mathbf{z} on \mathbf{X} .
3. Update the $\pi_1(x_i, \hat{\beta})$'s using (12.50) with the current estimate of $\hat{\beta}^{(k+1)}$.

If the data are *linearly separable*, that is, there exists a hyperplane that perfectly separates the two classes, the maximum conditional log-likelihood estimator for the logistic regression model does not exist (see Exercise 8). In this case, the iteratively reweighted least squares algorithm might not converge. On the other hand, the existence of the maximum conditional log-likelihood estimator is guaranteed if data are not linearly separable. In real applications, the data are rarely linear separable, so this would not be a problem.

We could also obtain the estimated standard errors of the final solution $\hat{\beta}$. For the k th iteration, recall that the Fisher information matrix $I(\hat{\beta}^{(k)})$ takes the form

$$I(\hat{\beta}^{(k)}) = -\mathbb{E} \left(\frac{\partial^2 \ell(\hat{\beta}^{(k)})}{\partial \beta \partial \beta^T} \right) \approx \mathbf{X}^T \mathbf{W} \mathbf{X}, \quad (12.65)$$

we estimate the standard error of $\hat{\beta}_j$ as the j th diagonal element of $I(\hat{\beta})^{-1}$.

12.9 Logistic Regression Versus LDA

There is a close connection between logistic regression and Gaussian LDA. Let (X, Y) be a pair of random variables where Y is binary and let $p_0(x) = p(x | Y = 0)$, $p_1(x) = p(x | Y = 1)$, $\pi_1 = \mathbb{P}(Y = 1)$. By Bayes' theorem,

$$\mathbb{P}(Y = 1 | X = x) = \frac{p(x | Y = 1)\pi_1}{p(x | Y = 1)\pi_1 + p(x | Y = 0)(1 - \pi_1)} \quad (12.66)$$

If we assume that each group is Gaussian with the same covariance matrix Σ , i.e., $X | Y = 0 \sim N(\mu_0, \Sigma)$ and $X | Y = 1 \sim N(\mu_1, \Sigma)$, we have

$$\log \left(\frac{\mathbb{P}(Y = 1 | X = x)}{\mathbb{P}(Y = 0 | X = x)} \right) = \log \left(\frac{\pi_1}{1 - \pi_1} \right) - \frac{1}{2}(\mu_0 + \mu_1)^T \Sigma^{-1}(\mu_1 - \mu_0) \quad (12.67)$$

$$+ x^T \Sigma^{-1}(\mu_1 - \mu_0) \quad (12.68)$$

$$\equiv \alpha_0 + \alpha^T x. \quad (12.69)$$

On the other hand, the logistic regression model is, by assumption,

$$\log \left(\frac{\mathbb{P}(Y = 1 | X = x)}{\mathbb{P}(Y = 0 | X = x)} \right) = \beta_0 + \beta^T x.$$

These are the same model since they both lead to classification rules that are linear in x . The difference is in how we estimate the parameters.

This is an example of a generative versus a discriminative model. In Gaussian LDA we estimate the whole joint distribution by maximizing the full likelihood

$$\prod_{i=1}^n p(x_i, y_i) = \underbrace{\prod_{i=1}^n p(x_i | y_i)}_{\text{Gaussian}} \underbrace{\prod_{i=1}^n p(y_i)}_{\text{Bernoulli}}. \quad (12.70)$$

In logistic regression we maximize the conditional likelihood $\prod_{i=1}^n p(y_i | X_i)$ but ignore the second term $p(x_i)$:

$$\prod_{i=1}^n p(x_i, y_i) = \underbrace{\prod_{i=1}^n p(y_i | x_i)}_{\text{logistic}} \underbrace{\prod_{i=1}^n p(x_i)}_{\text{ignored}}. \quad (12.71)$$

Since classification only requires the knowledge of $p(y|x)$, we don't really need to estimate the whole joint distribution. Logistic regression leaves the marginal distribution $p(x)$ unspecified so it relies on less parametric assumption than LDA. This is an advantage of the logistic regression approach over LDA. However, if the true class conditional distributions are Gaussian, logistic regression will be asymptotically less efficient than LDA, i.e., to achieve a certain level of classification error, logistic regression will require more data. In practice, logistic regression and LDA often give similar results.

12.10 Regularized Logistic Regression

As with linear regression, when the dimension d of the covariate is large, we cannot simply fit a logistic model to all the variables without experiencing numerical and statistical problems. Akin to the lasso, we will use *regularized logistic regression*, which includes *sparse logistic regression* and *ridge logistic regression*.

Let $\ell(\beta_0, \beta)$ be the log-likelihood defined in (12.54). The *sparse logistic regression* estimator is an ℓ_1 -regularized conditional log-likelihood estimator

$$\hat{\beta}_0, \hat{\beta} = \arg \min_{\beta_0, \beta} \left\{ -\ell(\beta_0, \beta) + \lambda \|\beta\|_1 \right\}. \quad (12.72)$$

Similarly, the *ridge logistic regression* estimator is an ℓ_2 -regularized conditional log-likelihood estimator

$$\hat{\beta}_0, \hat{\beta} = \arg \min_{\beta_0, \beta} \left\{ -\ell(\beta_0, \beta) + \lambda \|\beta\|_2^2 \right\}. \quad (12.73)$$

The algorithm for logistic ridge regression only requires a simple modification of the iteratively reweighted least squares algorithm and is left as an exercise (Exercise 9).

For sparse logistic regression, an easy way to calculate $\hat{\beta}_0$ and $\hat{\beta}$ is to apply a ℓ_1 -regularized Newton procedure. Similar to the Newton method for unregularized logistic regression, for the k th iteration, we first form a quadratic approximation to the negative log-likelihood $\ell(\beta_0, \beta)$ based on the current estimates $\hat{\beta}^{(k)}$.

$$-\ell(\beta_0, \beta) = \underbrace{\frac{1}{2} \sum_{i=1}^n w_{ii} (z_i^{(k)} - \beta_0 - \sum_{j=1}^d \beta_j x_{ij})^2}_{\ell_Q(\beta_0, \beta)} + \text{constant}. \quad (12.74)$$

where w_{ii} and $z_i^{(k)}$ are defined in (12.58) and (12.62). Since we have a ℓ_1 -regularization term, the updating formula for the estimate in the $(k+1)$ th step then becomes

$$\hat{\beta}^{(k+1)}, \hat{\beta}^{(k+1)} = \arg \min_{\beta_0, \beta} \left\{ \frac{1}{2} \sum_{i=1}^n w_{ii} (z_i^{(k)} - \beta_0 - \sum_{j=1}^d \beta_j x_{ij})^2 + \lambda \|\beta\|_1 \right\}. \quad (12.75)$$

This is a weighted lasso problem and can be solved using the coordinate descent algorithm we introduced before.

Sparse Logistic Regression Using Coordinate Descent

Choose starting values $\hat{\beta}^{(0)} = (\hat{\beta}_0^{(0)}, \hat{\beta}_1^{(0)}, \dots, \hat{\beta}_d^{(0)})^T$

(Outer loop) For $k = 1, 2, \dots$, iterate the following steps until convergence.

1. For $i = 1, \dots, n$, calculate $\pi_1(x_i, \hat{\beta}^{(k)})$, $z_i^{(k)}$, $w_{ii}^{(k)}$ according to (12.50), (12.62), and (12.58).

2. $\alpha_0 = \frac{\sum_{i=1}^n w_{ii}^{(k)} z_i^{(k)}}{\sum_{i=1}^n w_{ii}^{(k)}}$ and $\alpha_\ell = \hat{\beta}_\ell^{(k)}$ for $\ell = 1, \dots, d$.

3. (Inner loop) iterate the following steps until convergence

For $j \in \{1, \dots, d\}$

- (a) For $i = 1, \dots, n$, calculate $r_{ij} = z_i^{(k)} - \alpha_0 - \sum_{\ell \neq j} \alpha_\ell x_{i\ell}$.

- (b) Calculate $u_j^{(k)} = \sum_{i=1}^n w_{ii}^{(k)} r_{ij} x_{ij}$ and $v_j^{(k)} = \sum_{i=1}^n w_{ii}^{(k)} x_{ij}^2$.

- (c) $\alpha_j = \text{sign}(u_j^{(k)}) \left[\frac{|u_j^{(k)}| - \lambda}{v_j^{(k)}} \right]_+$.

4. $\hat{\beta}_0^{(k+1)} = \alpha_0$ and $\hat{\beta}_\ell^{(k+1)} = \alpha_\ell$ for $\ell = 1, \dots, d$.

5. Update the $\pi_1(x_i, \hat{\beta})$'s using (12.50) with the current estimate of $\hat{\beta}^{(k+1)}$.

Even though the above iterative procedure does not guarantee theoretical convergence, it works very well in practice.

12.11 Support Vector Machines

The *support vector machine* (SVM) classifier is like logistic regression, but it uses the *hinge loss* $L_{\text{hinge}}(y_i, H(x_i)) \equiv [1 - y_i H(x_i)]_+$ instead of the logistic loss. In this section, the outcomes are still coded as -1 and $+1$.

The support vector machine classifier is $\hat{h}(x) = I(\hat{H}(x) > 0)$ where the hyperplane $\hat{H}(x) = \hat{\beta}_0 + \hat{\beta}^T x$ is obtained by minimizing

$$\sum_{i=1}^n [1 - y_i H(x_i)]_+ + \frac{\lambda}{2} \|\beta\|_2^2 \quad (12.76)$$

where $\lambda > 0$ and the factor $1/2$ is only for notational convenience.

Figure 12.8 compares the hinge loss, 0-1 loss, and logistic loss. The advantage of the hinge loss is that it is convex, and it has a corner which leads to efficient computation and the minimizer of $\mathbb{E}(1 - YH(X))_+$ is the Bayes rule.¹ A disadvantage of the hinge loss is that one can't recover the regression function $m(x) = \mathbb{E}(Y | X = x)$. This is because $\hat{H}(x)$ is estimating $\text{sign}(m(x) - 1/2)$ not $m(x)$. To quote Hastie and Zhu (2006): “The SVM is not fundamentally different from many statistical tools that statisticians are familiar with, for example, penalized logistic regression.”

The SVM classifier is traditionally developed from a geometric perspective. Suppose first that the data are *linearly separable*, that is, there exists a hyperplane that perfectly separates the two classes. How can we find a separating hyperplane? LDA is not guaranteed to find it (See Exercise 13). A separating hyperplane will minimize

$$-\sum_{i \in \mathcal{M}} y_i H(x_i).$$

where \mathcal{M} is the index set of all misclassified data points. Rosenblatt's perceptron algorithm² takes starting values and iteratively updates the coefficients as:

$$\begin{pmatrix} \beta \\ \beta_0 \end{pmatrix} \leftarrow \begin{pmatrix} \beta \\ \beta_0 \end{pmatrix} + \rho \begin{pmatrix} y_i x_i \\ y_i \end{pmatrix}$$

where $\rho > 0$ is the learning rate. If the data are linearly separable, the perceptron algorithm is guaranteed to converge to a separating hyperplane. However, there could be many separating hyperplanes. Different starting values may lead to different separating hyperplanes. The question is, which separating hyperplane is the best?

Intuitively, it seems reasonable to choose the hyperplane “furthest” from the data in the sense that it separates the +1's and -1's and maximizes the distance to the closest point. This hyperplane is called the *maximum margin hyperplane*. The margin is the distance from the hyperplane to the nearest data point³. Points on the boundary of the margin are called *support vectors*. See Figure 12.9. The goal, then, is to find a separating hyperplane which maximizes the margin. After some simple algebra, we can show that (12.76) exactly achieves this goal. In fact, (12.76) also works for data that are not linearly separable. A detailed discussion will be provided in the later optimization section.

¹ The hinge loss is an example of a surrogate loss. We will discuss the theoretical properties of surrogate loss functions in detail later.

² more details about the perceptron algorithm will be discussed in the online learning section.

³ Some authors also define the margin to be twice the distance from the hyperplane to the nearest data point.

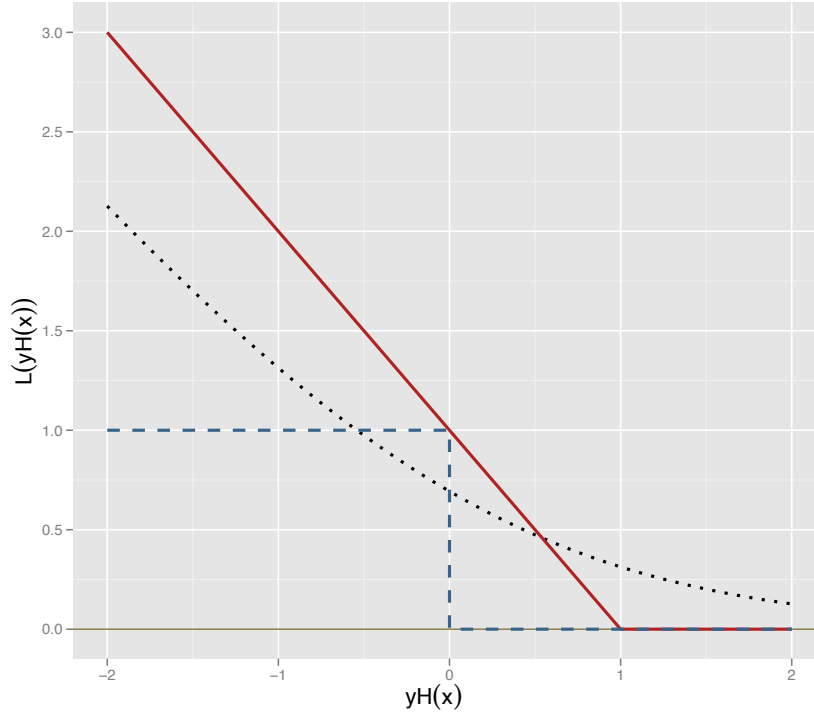


Figure 12.8. The 0-1 classification loss (blue dashed line), hinge loss (red solid line) and logistic loss (black dotted line).

The unconstrained optimization problem (12.76) can be equivalently formulated in constrained form:

$$\min_{\beta_0, \beta} \left\{ \frac{1}{2} \|\beta\|_2^2 + \frac{1}{\lambda} \sum_{i=1}^n \xi_i \right\} \quad (12.77)$$

$$\text{subject to } \forall i, \xi_i \geq 0 \text{ and } \xi_i \geq 1 - y_i H(x_i). \quad (12.78)$$

Given two vectors a and b , let $\langle a, b \rangle = a^T b = \sum_j a_j b_j$ denote the inner product of a and b . The following lemma provides the dual of the optimization problem in (12.77).

12.79 Lemma. *The dual of the SVM optimization problem in (12.77) takes the form*

$$\hat{\alpha} = \arg \max_{\alpha \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \langle x_i, x_k \rangle \right\} \quad (12.80)$$

$$\text{subject to } 0 \leq \alpha_1, \dots, \alpha_n \leq \frac{1}{\lambda} \text{ and } \sum_{i=1}^n \alpha_i y_i = 0, \quad (12.81)$$

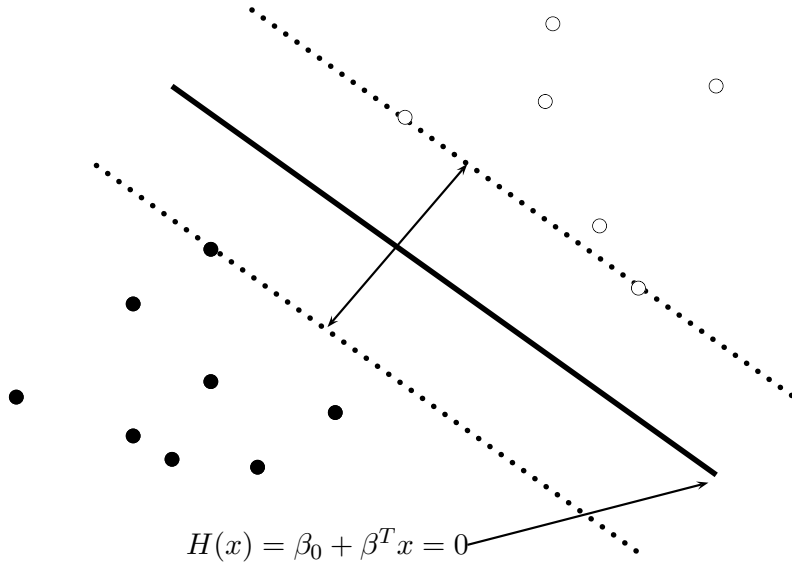


Figure 12.9. The hyperplane $H(x)$ has the largest margin of all hyperplanes that separate the two classes.

with the primal-dual relationship $\hat{\beta} = \sum_{i=1}^n \hat{\alpha} y_i x_i$. We also have

$$\hat{\alpha}_i (1 - \xi_i - y_i (\hat{\beta}_0 + \hat{\beta} x_i)) = 0, \quad i = 1, \dots, n. \quad (12.82)$$

Proof. Let $\alpha_i, \gamma_i \geq 0$ be the Lagrange multipliers. The Lagrangian function can be written as

$$L(\xi, \beta, \beta_0, \alpha, \gamma) = \frac{1}{2} \|\beta\|_2^2 + \frac{1}{\lambda} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - \xi_i - y_i H(x_i)) - \sum_{i=1}^n \gamma_i \xi_i. \quad (12.83)$$

The Karush-Kuhn-Tucker conditions are

$$\forall i, \alpha_i \geq 0, \gamma_i \geq 0, \xi_i \geq 0 \text{ and } \xi_i \geq 1 - y_i H(x_i), \quad (12.84)$$

$$\beta = \sum_{i=1}^n \alpha_i y_i x_i, \quad \sum_{i=1}^n \alpha_i y_i = 0 \text{ and } \alpha_i + \gamma_i = 1/\lambda, \quad (12.85)$$

$$\forall i, \alpha_i (1 - \xi_i - y_i H(x_i)) = 0 \text{ and } \gamma_i \xi_i = 0. \quad (12.86)$$

The dual formulation in (12.80) follows by plugging (12.84) and (12.85) into (12.83). The primal-dual complementary slackness condition (12.82) is obtained from the first equation in (12.86). \square

The dual problem (12.80) is easier to solve than the primal problem (12.76). The data points (x_i, y_i) for which $\hat{\alpha}_i > 0$ are called *support vectors*. By (12.82) and (12.77), for all the data points (x_i, y_i) satisfying $y_i(\hat{\beta}_0 + \hat{\beta}^T x_i) > 1$, there must be $\hat{\alpha}_i = 0$. The solution for the dual problem is sparse. From the first equality in (12.85), we see that the final estimate $\hat{\beta}$ is a linear combination only of these support vectors. Among these support vectors, if $\alpha_i < 1/\lambda$, we call (x_i, y_i) a *margin point*. For a margin point (x_i, y_i) , the last equality in (12.85) implies that $\gamma_i > 0$, then the second equality in (12.86) implies $\xi_i = 0$. Moreover, using the first equality in (12.86), we get

$$\hat{\beta}_0 = -y_i x_i^T \hat{\beta}. \quad (12.87)$$

Therefore, once $\hat{\beta}$ is given, we could calculate $\hat{\beta}_0$ using any margin point (x_i, y_i) .

12.88 Example. We consider classifying two types of irises, versicolor and virginica. There are 50 observations in each class. The covariates are "Sepal.Length" "Sepal.Width" "Petal.Length" and "Petal.Width". After fitting a SVM we get a 3/100 misclassification rate. The SVM uses 33 support vectors. \square

12.12 Fitting Support Vector Machines

The sequential minimal optimization (SMO) algorithm, due to Platt (1998), provides an efficient solution to the the dual form SVM problem in (12.80). Here, we introduce a slight modification of this algorithm due to Bottou and Lin (2007).

Let $C = 1/\lambda$. The dual form SVM problem can be written as

$$\max_{\alpha} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^n \alpha_i \alpha_k y_i y_k \langle x_i, x_k \rangle \right\} \quad (12.89)$$

$$\text{subject to } 0 \leq \alpha_i \leq C \text{ for } i = 1, \dots, n, \quad (12.90)$$

$$\sum_{i=1}^n \alpha_i y_i = 0. \quad (12.91)$$

The SMO algorithm adopts the *coordinate ascent* strategy. In each iteration, we hold all except a subset of α 's fixed, and optimize the objective function in (12.89) with respect to just the variables in this subset. In general, the coordinate ascent algorithm updates one variable at a time. However, due to the equality constraint in (12.91), any variable α_i will be uniquely determined by the remaining $n - 1$ variables since $\alpha_i = -y_i \sum_{j \neq i} y_j \alpha_j$. Therefore, the SMO algorithm updates two variables within each iteration. This adds some complications. Compared to the classical coordinate ascent method, we need to answer two questions: (i) which pair of variables to update within each iteration? (ii) how to design a criterion to test the convergence of the algorithm? Different answers to these two questions result in different variants of the SMO algorithm. Here we introduce a simple version due to Bottou and Lin (2007).

We define $L_i = -C \cdot I(y_i = -1)$ and $U_i = C \cdot I(y_i = +1)$. It is obvious that $y_i \alpha_i \in [L_i, U_i]$ for $i = 1, \dots, n$. The SMO algorithm is described in the following figure. A detailed derivation of this algorithm can be found in (Bottou and Lin, 2007).

SVM Using Sequential Minimal Optimization Algorithm

Initialize coefficients and gradients, for all $i = 1, \dots, n$, $\alpha_i = 0$ and $g_i = 0$.

Iterate the following steps until step 3 is satisfied or the maximum iteration number has been reached:

1. $i = \arg \min_i y_i g_i$ subject to $y_i \alpha_i < U_i$.
2. $j = \arg \min_j y_j g_j$ subject to $L_j < y_j \alpha_j$.
3. If $y_i g_i \leq y_j g_j$ then stop.
4. $\rho = \min \left\{ U_i - y_i \alpha_i, y_j \alpha_j - L_j, \frac{y_i g_i - y_j g_j}{\langle x_i, x_i \rangle + \langle x_j, x_j \rangle - 2 \langle x_i, x_j \rangle} \right\}$.
5. For all $k = 1, \dots, n$, $g_k \leftarrow g_k - \rho y_k \langle x_i, x_k \rangle + \rho y_k \langle x_j, x_k \rangle$.
6. $\alpha_i \leftarrow \alpha_i + \rho y_i$ and $\alpha_j \leftarrow \alpha_j - \rho y_j$.

As an alternative to the above SMO algorithm which solves the dual problem (12.89), Shalev-Shwartz et al. (2007) develop a stochastic subgradient projection method to directly solve the unconstrained SVM optimization (12.76). Their algorithm, called Pegasos, is simple and scalable to very large datasets.

12.13 Case Study I: Supernova Classification

A *supernova* is an exploding star. Type Ia supernovae are a special class of supernovae that are very useful in astrophysics research. These supernovae have a characteristic *light curve*, which is a plot of the luminosity of the supernova versus time. The maximum brightness of all type Ia supernovae is approximately the same. In other words, the true (or absolute) brightness of a type Ia supernova is known. On the other hand, the apparent (or observed) brightness of a supernova can be measured directly. Since we know both the absolute and apparent brightness of a type Ia supernova, we can compute its distance. Because of this, type Ia supernovae are sometimes called standard candles. Two supernovae, one type Ia and one non-type Ia, are illustrated in Figure 12.10. Astronomers also measure the *redshift* of the supernova, which is essentially the speed at which the supernova is moving away from us. The relationship between distance and redshift provides important information for astrophysicists in studying the large scale structure of the universe.

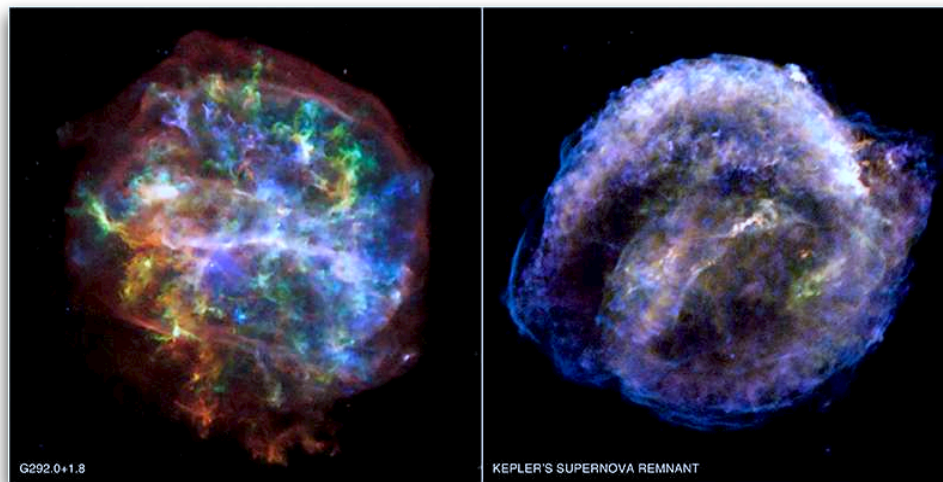


Figure 12.10. Two supernova remnants from the NASA's Chandra X-ray Observatory study. The image in the right panel, the so-called Kepler supernova remnant, is "Type Ia". Such supernovae have a very symmetric, circular remnant. This type of supernova is thought to be caused by a thermonuclear explosion of a white dwarf, and is often used by astronomers as a "standard candle" for measuring cosmic distances. On the other hand, the image in the left panel is the morphology of the G292.0+1.8 remnant, comes from the "core collapse" family of supernova explosions. Such supernovae are distinctly more asymmetric. (Credit: NASA/CXC/UCSC/L. Lopez et al.)

A major challenge in astrophysics is to classify supernovae to be type Ia versus other types. Kessler et al. (2010) released a mixture of real and realistically simulated supernovae and challenged the scientific community to find effective ways to classify the type Ia supernovae. The dataset consists of about 20,000 simulated supernovae. For each supernova, there are a few noisy measurements of the flux (brightness) in four different filters. These four filters correspond to different wavelengths. Specifically, the filters correspond to the g -band (green), r -band (red), i -band (infrared) and z -band (blue). See Figure 12.11.

To estimate a linear classifier we need to preprocess the data to extract features. One difficulty is that each supernova is only measured at a few irregular time points, and these time points are not aligned. To handle this problem, using the estimated measurement errors of each flux as weights, we estimate a weighted least squares regression spline (see later chapters on nonparametric methods) to smooth each supernova. All four filters of each supernova are then aligned according to the peak of the r -band. We also rescale so that all the curves have the same maximum.

The goal of this study is to build linear classifiers to predict whether a supernova is type Ia or not. For simplicity, we only use the information in the r -band. First, we align the fitted regression spline curves of all supernovae by calibrating their maximum peaks and set the corresponding time point to be day 0. There are altogether 19,679 supernovae in the dataset with 1,367 being labeled. To get a higher signal-to-noise ratio, we throw away all supernovae with less than 10 r -band flux measurements. We finally get a trimmed dataset

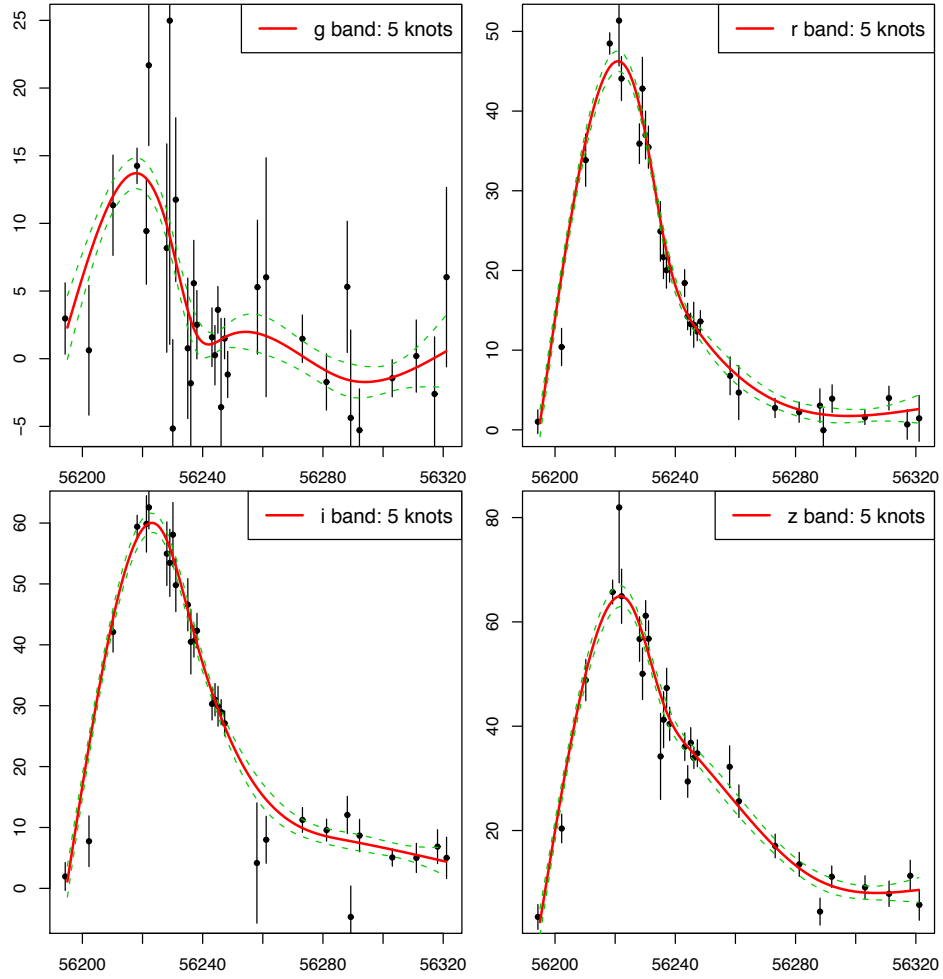


Figure 12.11. Four filters (g , r , i , z -bands) corresponding to a type Ia supernova DES-SN000051. For each band, a weighted regression spline fit (solid red) with the corresponding standard error curves (dashed green) is provided. The black points with bars represent the flux values and their estimated standard errors.

with 255 supernovae, 206 of which are type Ia and 49 of which are non-type Ia.

We use two types of features: the *time-domain* features and *frequency-domain* features. For the time-domain features, the features are the interpolated regression spline values according to an equally spaced time grid. In this study, the grid has length 100, ranging from day -20 to day 80. Since all the fitted regression curves have similar global shapes, the time-domain features are expected to be highly correlated. This conjecture is confirmed by the scatter matrix of the first five features in 12.12. To make the features less correlated, we also extract the frequency-domain features, which are simply the discrete cosine trans-

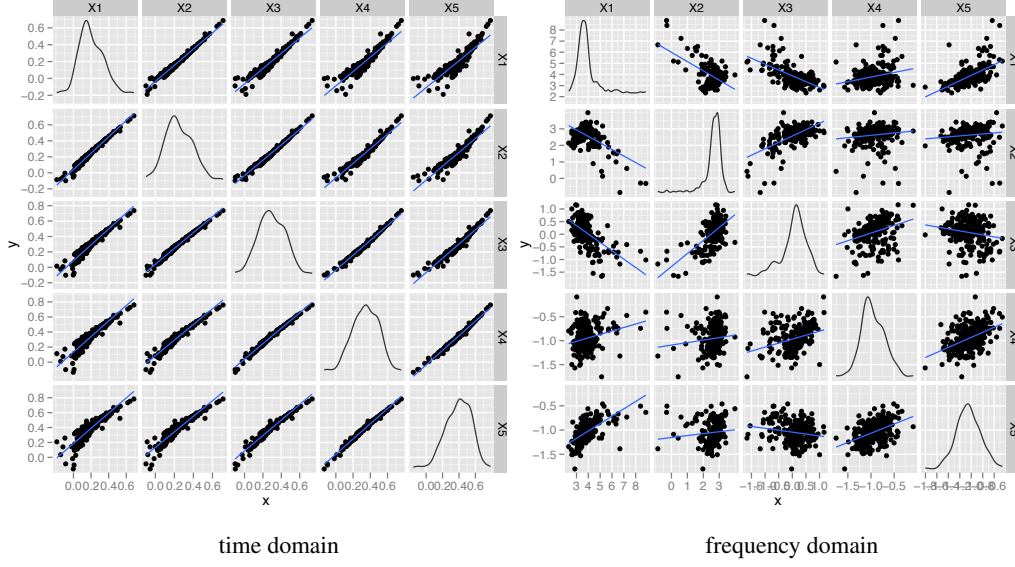


Figure 12.12. The scatter matrix of the first five features of the supernova data. On the diagonal cells are the estimated univariate densities of each feature. The off-diagonal cells visualize the pairwise scatter plots of the two corresponding variables with a least squares fit. We see the time-domain features are highly correlated, while the frequency-domain features are almost uncorrelated.

formations of the corresponding time-domain features. More specifically, given the time domain features X_1, \dots, X_d ($d = 100$), Their corresponding frequency domain features $\tilde{X}_1, \dots, \tilde{X}_d$ can be written as

$$\tilde{X}_j = \frac{2}{d} \sum_{k=1}^d X_k \cos \left[\frac{\pi}{d} \left(k - \frac{1}{2} \right) (j - 1) \right] \text{ for } j = 1, \dots, d. \quad (12.92)$$

The right panel of Figure 12.12 illustrates the scatter matrix of the first 5 frequency-domain features. In contrast to the time-domain features, the frequency-domain features have low correlation.

We apply sparse logistic regression (LR), support vector machines (SVM), diagonal linear discriminant analysis (DLDA), and diagonal quadratic discriminant analysis (DQDA) on this dataset. For each method, we conduct 100 runs, within each run, 40% of the data are randomly selected as training and the remaining 60% are used for testing.

Figure 12.13 illustrates the regularization paths of sparse logistic regression using the time-domain and frequency-domain features. A regularization path provides the coefficient value of each feature over all regularization parameters. Since the time-domain features are highly correlated, the corresponding regularization path is quite irregular. In contrast, the paths for the frequency-domain features behave stably.

Figure 12.14 compares the classification performance of all these methods. The results show that classification in the frequency domain is not helpful. The regularization paths of the SVM are the same in both the time and frequency domains. This is expectable since

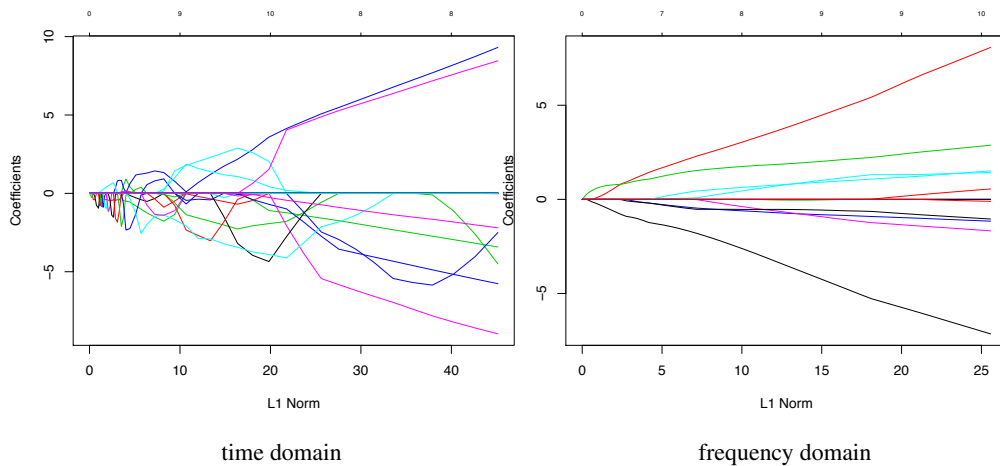


Figure 12.13. The regularization paths of sparse logistic regression using the features of time-domain and frequency-domain. The vertical axis corresponds to the values of the coefficients, plotted as a function of their ℓ_1 -norm. The path using time-domain features are highly irregular, while the path using frequency-domain features are more stable.

the discrete cosine transformation is an orthonormal transformation, which corresponds to rotate the data in the feature space while preserving their Euclidean distances and inner products. From (12.89), it is easy to see that the SVM is rotation invariant. The sparse logistic regression is not rotation invariant due to the ℓ_1 -norm regularization term. The performance of the sparse logistic regression in the frequency domain is worse than that in the time domain. The DLDA and DQDA are also not rotation invariant, their performances decrease significantly in the frequency domain than those in the time domain. In both time and frequency domains, the SVM outperforms all the other methods. Then follows sparse logistic regression, which is better than DLDA and DQDA. At the first sight, it might be counterintuitive that, using the time-domain features, the DQDA has a higher training error than the DLDA, recalling that DQDA has more free parameters than DLDA. However, on the training sets, the DQDA tries to maximize the Gaussian likelihood, which is not directly related to minimizing the classification error.

12.14 Case Study II: Political Blog Classification

A *blog* is a special type of website that is usually maintained by an individual who regularly posts entries of commentary or descriptions of events. Blogs are made for public access and their entries are usually displayed in reverse-chronological order. In this example, we classify political blogs according to whether their political leanings are *liberal* or *conservative*. Snapshots of two political blogs are shown in Figure 12.15.

A corpus of 403 political blogs is collected by Lin and Cohen (2010) for a two-month window before the 2004 presidential election. Among these blogs, 205 are liberal and 198 are conservative. We use *bag-of-words* features, i.e., each unique word from these 403 blogs

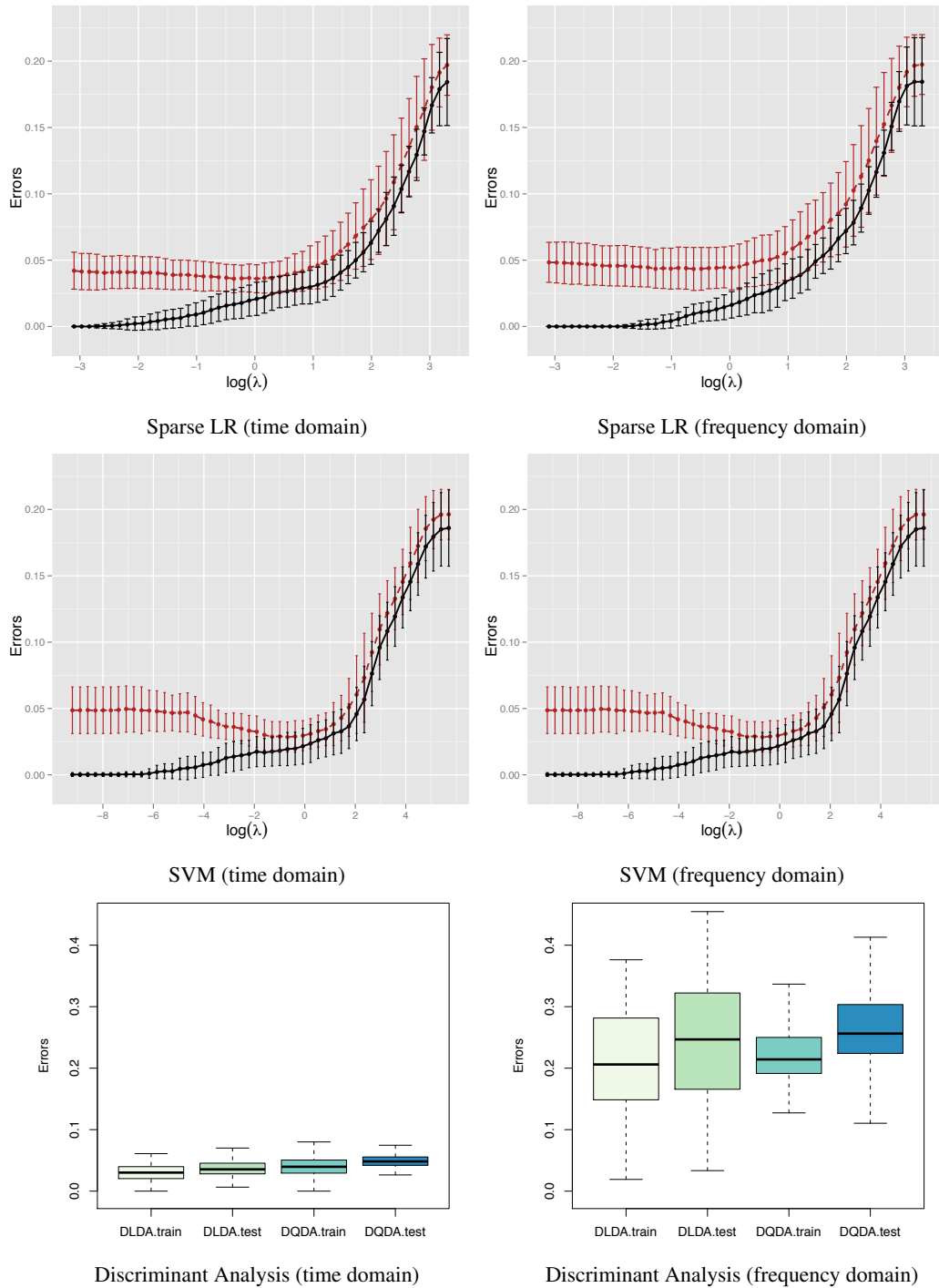


Figure 12.14. Comparison of different methods on the supernova dataset using both the time-domain (left column) and frequency-domain features (right Column). Top four figures: mean error curves (black: training error; red: test error) and their corresponding standard error bars for sparse logistic regression (LR) and support vector machines (SVM). Bottom two figures: boxplots of the training and test errors of diagonal linear discriminant analysis (DLDA) and diagonal quadratic discriminant analysis (DQDA). For the time-domain features, the SVM achieves the smallest test error among all methods.

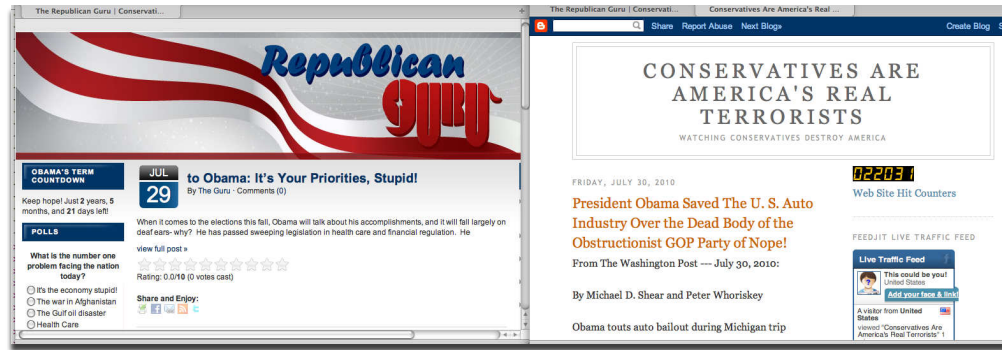


Figure 12.15. Examples of two political blogs with different orientations, one conservative and the other liberal.

serves as a feature. For each blog, the value of a feature is the number of occurrences of the word normalized by the total number of words in the blog. After converting all words to lower case, we remove stop words and only retain words with at least 10 occurrences across all the 403 blogs. This results in 23,955 features, each of which corresponds to an English word. Such features are only a crude representation of the text represented as an unordered collection of words, disregarding all grammatical structure. We also extracted features that use hyperlink information. In particular, we selected 292 out of the 403 blogs that are heavily linked to, and for each blog $i = 1, \dots, 403$, its linkage information is represented as a 292-dimensional binary vector $(x_{i1}, \dots, x_{i292})^T$ where $x_{ij} = 1$ if the i th blog has a link to the j th feature blog. The total number of covariates is then $23,955 + 292 = 24,247$. Even though the link features only constitute a small proportion, they are important for predictive accuracy.

We run the full regularization paths of sparse logistic regression and support vector machines, 100 times each. For each run, the data are randomly partitioned into training (60%) and testing (40%) sets. Figure 12.16 shows the mean error curves with their standard errors. From Figure 12.16, we see that linkage information is crucial. Without the linkage features, the smallest mean test error of the support vector machine along the regularization path is 0.247, while that of the sparse logistic regression is 0.270. With the link features, the smallest test error for the support vector machine becomes 0.132. Although the support vector machine has a better mean error curve, it has much larger standard error. Two typical regularization paths for sparse logistic regression with and without using the link features are provided at the bottom of Figure 12.16. By examining these paths, we see that when the link features are used, 11 of the first 20 selected features are link features. In this case, although the class conditional distribution is obviously not Gaussian, we still apply the diagonal linear discriminant analysis (DLDA) on this dataset for a comparative study. Without the linkage features, the DLDA has a mean test error 0.303 (sd = 0.07). With the linkage features, DLDA has a mean test error 0.159 (sd = 0.02).

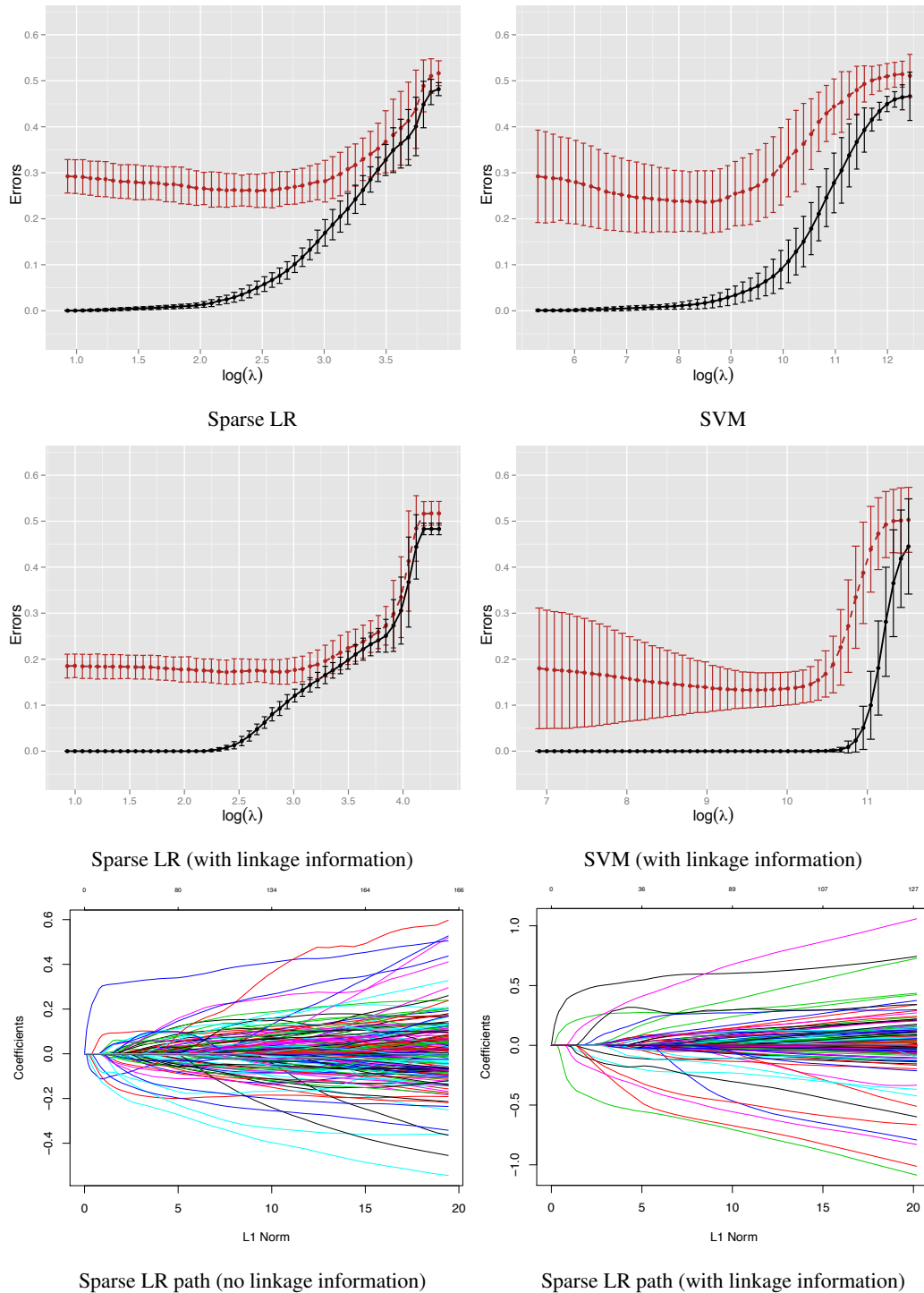


Figure 12.16. Comparison of the sparse logistic regression (LR) and support vector machine (SVM) on the political blog data. (Top four figures): The mean error curves (Black: training error, Red: test error) and their corresponding standard error bars of the sparse LR and SVM, with and without the linkage information. (Bottom two figures): Two typical regularization paths of the sparse logistic LR with and without the linkage information. On this dataset, the diagonal linear discriminant analysis (DLDA) achieves a test error 0.303 ($sd = 0.07$) without the linkage information and a test error 0.159 ($sd = 0.02$) with the linkage information.

12.15 Bibliographic Remarks

References on linear discriminant analysis include Devroye et al. (1996) and Duda et al. (2000). References on logistic regression include Agresti (1990), Dobson (2001), and McCullagh and Nelder (1999). Recent progress on sparse classification methods is described in Hastie et al. (2009).

The theory of SVMs is developed in Vapnik (1995) and Vapnik (1998). Extension of the SVM to regression settings can be found in Smola and Schölkopf (2004). A quick summary of multiclass SVMs is provided in Liu (2007).

Exercises

- 12.1 Suppose that $\mathbb{P}(Y = 1) = \mathbb{P}(Y = 0) = \frac{1}{2}$ and $X | Y = 0 \sim N(0, 1)$ and $X | Y = 1 \sim \frac{1}{2}N(-5, 1) + \frac{1}{2}N(5, 1)$.
- Find expressions for the Bayes classifier and the Bayes risk.
 - What linear classifier minimizes the risk and what is its risk?
- 12.2 Suppose that $\mathbb{P}(Y = 1) = \mathbb{P}(Y = -1) = \frac{1}{2}$ and $X | Y = -1 \sim \text{Uniform}(-10, 5)$ and $X | Y = 1 \sim \text{Uniform}(-5, 10)$.
- Find an expression for the Bayes classifier and find an expression for the Bayes risk.
 - Consider the linear classifier $h_\beta(x) = \text{sign}(\beta x)$ where $\beta \in \mathbb{R}$. What linear classifier β^* minimizes the risk and what is its risk?
 - Compute the hinge risk $R_\phi(\beta) = \mathbb{E}(1 - Y\beta X)_+$.
- 12.3 We define that $X \in \mathbb{R}$ has a *nonparanormal* distribution, written as $X \sim \text{NPN}(\mu, \sigma^2, f)$, if f is a monotonic increasing function and $f(X) \sim N(\mu, \sigma^2)$. Suppose that $\mathbb{P}(Y = 1) = \mathbb{P}(Y = 0) = \frac{1}{2}$ and
- $$X | Y = 0 \sim \text{NPN}(0, 1, f)$$
- $$X | Y = 1 \sim \text{NPN}(1, 1, f).$$
- Find an expression for the Bayes classifier and find an expression for the Bayes risk.
 - What linear classifier minimizes the risk and what is its risk?
 - Give an algorithm to estimate the best linear classifier from a sample $(x_1, y_1), \dots, (x_n, y_n)$.
- 12.4 Prove Equations (12.27) to (12.30).
- 12.5 Prove that Equation (12.33) is the maximum likelihood estimate of the common covariance matrix Σ .

- 12.6 Prove Theorem 12.40.
- 12.7 Show that if Newton's method is applied to the logistic regression log-likelihood, it leads to the reweighted least squares algorithm.
- 12.8 Show that if the data are perfectly separable, the conditional maximum likelihood estimator for the logistic regression model does not exist. Comment on the behavior of the iteratively reweighted least squares algorithm.
- 12.9 Derive the Newton algorithm for ridge logistic regression in (12.73). Compare it with the iteratively reweighted least squares algorithm and comment the difference.
- 12.10 In this problem you are asked to fit a logistic regression model to the UCI Pima Indians diabetes database. The data, and a description of the data, can be downloaded from <http://ftp.ics.uci.edu/pub/machine-learning-databases/pima-indians-diabetes/>. It is a binary classification problem, with 768 instances having eight features each.
- (a) Fit a maximum likelihood logistic regression model using Newton's method (iteratively reweighted least squares). Train the model on subsets of size $L = 10, 20, 30, \dots, 100$ examples. Plot a curve of accuracy versus L , averaging over 20 random choices of the training set for each L (and testing on the remainder). Show the standard errors in your plot.
 - (b) For $L = 100$, give a plot of training and test set log-likelihood as a function of Newton iteration.
- 12.11 Get the Coronary Risk-Factor Study (CORIS) data from the book web site. Use backward stepwise logistic regression based on AIC to select a model. Summarize your results.
- 12.12 Derive the coordinate descent algorithm for sparse logistic regression.
- 12.13 Construct a concrete binary class classification example, in which the data from the two classes are linear separable but the LDA solution does not separate the data.
- 12.14 Suppose $Y \in \{0, 1\}$ is a random variable given by

$$Y = \begin{cases} 1 & a^\top u + b + v \leq 0 \\ 0 & a^\top u + b + v > 0 \end{cases}$$

where $u \in \mathbb{R}^n$ is a vector of explanatory variables and v is distributed as a zero mean unit variance Gaussian variable. Formulate the maximum likelihood estimation problem of estimating a and b , given data consisting of pairs (u_i, y_i) , $i = 1, \dots, n$, as a convex optimization problem.

- 12.15 Linear regression models a real-valued output Y given an input vector X as

$$Y | X \sim \text{Normal}(\mu(X), \sigma^2)$$

where the mean is a linear function of the input:

$$\mu(X) = \beta^T X = \beta_0 + \beta_1 X_1 + \dots + \beta_d X_d.$$

In logistic regression, the conditional distribution of a binary output Y is Bernoulli:

$$Y | X \sim \text{Binom}(\theta(X))$$

where the Bernoulli parameter is related to $\beta^T X$ by the logit transformation

$$\text{logit}(\theta(X)) \equiv \log \left(\frac{\theta(X)}{1 - \theta(X)} \right) = \beta^T X$$

If the output Y is a “count,” which can take any positive integer value $0, 1, 2, \dots$ (for example, the number of daily hits on a web server), the standard model is the Poisson. Recall that the probability mass function of the Poisson with parameter $\lambda > 0$ is

$$P(Y = k | \lambda) = \frac{\lambda^k}{k!} e^{-\lambda}$$

In *Poisson regression* the conditional distribution of Y given X is

$$Y | X \sim \text{Poisson}(\lambda(X))$$

where

$$\log \lambda(X) = \beta^T X$$

- (a) For data $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$ where $x_i \in \mathbb{R}^p$ and $y_i \in \mathbb{N} = \{0, 1, 2, \dots\}$, write down the (conditional) log-likelihood function $\ell(\beta)$ under the above Poisson regression model.
- (b) For each of linear regression, logistic regression, and Poisson regression, show that at the MLE $\hat{\beta}$

$$\sum_{i=1}^n y_i x_i = \sum_{i=1}^n E_{\hat{\beta}}[Y | X = x_i] x_i$$

- (c) Give an algorithm for computing the MLE $\hat{\beta}$ for Poisson regression. Your algorithm should be explicit: either give a closed form expression for the MLE as a function of the data, or give an iterative algorithm by specifying the update for $\beta^{(k)}$ in terms of the data and $\beta^{(k-1)}$.

Bibliography

- AGRESTI, A. (1990). *Categorical Data Analysis*. Wiley.
- BICKEL, P. J. and LEVINA, E. (2004). Some theory for fisher's linear discriminant function, 'naive bayes', and some alternatives when there are many more variables than observations. *Bernoulli* **10** 989–1010.
- BOTTOU, L. and LIN, C.-J. (2007). Support vector machine solvers. In *Large Scale Kernel Machines* (L. Bottou, O. Chapelle, D. DeCoste and J. Weston, eds.). MIT Press, Cambridge, MA., 301–320.
URL <http://leon.bottou.org/papers/bottou-lin-2006>
- DEVROYE, L., GYÖRFI, L. and LUGOSI, G. (1996). *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag. New York, NY.
- DOBSON, A. J. (2001). *An introduction to generalized linear models*. Chapman & Hall.
- DUDA, R. O., HART, P. E. and STORK, D. G. (2000). *Pattern Classification (2nd Edition)*. 2nd ed. Wiley-Interscience.
- HASTIE, T., TIBSHIRANI, R. and FRIEDMAN, J. H. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction (Second Edition)*. Springer-Verlag.
- KESSLER, R., CONLEY, A., JHA, S. and KUHLMANN, S. (2010). Supernova photometric classification challenge. Tech. Rep. arXiv:1001.5210.
- LIN, F. and COHEN, W. (2010). Power iteration clustering. In *Proceedings of ICML-10, 27th International Conference on Machine Learning*.
- LIU, Y. (2007). Fisher consistency of multicategory support vector machines. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, vol. 11.
- MCCULLAGH, P. and NELDER, J. A. (1999). *Generalized linear models*. Chapman and Hall. New York, NY.
- PLATT, J. C. (1998). Fast training of support vector machines using sequential minimal optimization. In *Advances in Kernel Methods: Support Vector Learning*. MIT Press, Cambridge, MA.

- SHALEV-SHWARTZ, S., SINGER, Y. and SREBRO, N. (2007). Pegasos: Primal estimated sub-gradient solver for svm. In *Proceedings of the Twenty-Fourth International Conference*.
- SMOLA, A. and SCHÖLKOPF, B. (2004). A tutorial on support vector regression. *Statistics and Computing* **14** 199–222.
- VAPNIK, V. (1995). *The Nature of Statistical Learning Theory*. 2nd ed. Springer.
- VAPNIK, V. N. (1998). *Statistical learning theory*. Wiley.