



Published in final edited form as:

*J Am Stat Assoc.* 2012 September 1; 107(499): 1166–1177. doi:10.1080/01621459.2012.699793.

## Deconvolution When Classifying Noisy Data Involving Transformations

**Raymond Carroll [Head],**

Department of Statistics, Texas A&M University, College Station, TX 77843-3143

**Aurore Delaigle [Associate Professor], and**

Department of Mathematics and Statistics, University of Melbourne, VIC 3010, Australia

**Peter Hall [Professor]**

Department of Mathematics and Statistics, University of Melbourne, VIC 3010, Australia

Raymond Carroll: carroll@stat.tamu.edu; Aurore Delaigle: a.delaigle@ms.unimelb.edu.au; Peter Hall: halpstat@ms.unimelb.edu.au

### Abstract

In the present study, we consider the problem of classifying spatial data distorted by a linear transformation or convolution and contaminated by additive random noise. In this setting, we show that classifier performance can be improved if we carefully invert the data before the classifier is applied. However, the inverse transformation is not constructed so as to recover the original signal, and in fact, we show that taking the latter approach is generally inadvisable. We introduce a fully data-driven procedure based on cross-validation, and use several classifiers to illustrate numerical properties of our approach. Theoretical arguments are given in support of our claims. Our procedure is applied to data generated by light detection and ranging (Lidar) technology, where we improve on earlier approaches to classifying aerosols. This article has supplementary materials online.

### Keywords

Centroid classifier; Cross-validation; Fourier transform; Inverse transform; Spatial data

## 1. INTRODUCTION

In the present study, we consider signal classification problems, where the observations are  $d$ -dimensional noisy spatial functions  $Y_{ij}$ , for  $1 \leq i \leq n_j$ , coming from population  $\Pi_j$ , where  $j = 1$  or  $2$ , and which can be modeled as  $Y_{ij} = TX_{ij} + \delta_{ij}$ , where  $T$  is a transformation of the function of interest  $X_{ij}$  and  $\delta_{ij}$  is a random error with zero mean and some correlation structure. Based on training data, the goal is to classify a new noisy data function  $Y$  whose class is unknown, as coming from one of  $\Pi_1$  and  $\Pi_2$ .

In many instances, the function  $TX_{ij}$  is the result of a convolution of the function  $X_{ij}$  with a blurring source, that is,  $TX_{ij} = \omega_T * X_{ij}$  where  $*$  denotes the convolution operator (see Section 2.3) and  $\omega_T$  is a point spread function. There, the function  $X_{ij}$  can be reconstructed in part by a (necessarily estimated) deconvolution operation. There is a large statistics

literature on deconvolution for image data, and for data of similar type, dating from the 1980s. It includes contributions by Besag (1986), Donoho (1994), Dass and Nair (2003), Qiu (2005, 2007, 2008), and Mukherjee and Qiu (2011). Related research problems arise in spatial statistics, for example, in the contexts of remote sensing (see, e.g., Klein and Press 1992; Cressie and Kornak 2003; Crosilla, Visintini, and Sepic 2007) and statistical signal recovery (see, e.g., Johnstone 1990; Huang and Cressie 2000; Shi and Cressie 2007).

There is also a significant literature on blind deconvolution and estimation of point spread functions. This work includes contributions by Kundur and Hatzinakos (1998), Cannon (1976), Carasso (2001), Galatsanos et al. (2002), Figueiredo and Nowak (2003), Joshi and Chaudhuri (2005), Hall and Qiu (2007a, b), Qiu (2008), Huang and Qiu (2010), and Popescu and Hellicar (2010). However, the problems of deconvolution and point spread function estimation are very different from those of classification, to such an extent that, even if  $T$  were known, the methods suggested in this article would still be recommended. It should also be noted that, since neither the function  $X$  nor the noise  $\delta$  is observable, it is not possible to estimate the noise and, hence, to remove it effectively from the observed data  $Y$ . In particular, in the problem treated in this article, it is not possible to compute residuals.

In our classification context, it is at least intuitively plausible that if one could recover the function  $X_{ij}$  then one would use that function as the basis for classification, rather than using the noisy convolved function  $Y_{ij}$ . This idea has been used in the classification of different types of aerosols using long-range infrared light detection and ranging (Lidar) methods (Warren et al. 2008), where deconvolution was used to obtain estimates of the true signal, and the resulting estimates were used as the basis for classification. Our work relates to whether a signal should be deconvolved or correlated errors should be deconvolved before classification, and we shall use Lidar data to illustrate our conclusions. We shall show that there exists a transformation of the noisy convolved function  $Y_{ij}$  that is appropriate for classification, but that it is not necessarily related to the transformation that would be used to recover the true signal.

The real-data classification problems that motivate this work, all involve just  $K = 2$  populations, and for this reason, and to simplify discussion, we shall confine attention to this case. However, our methodology and theoretical results extend readily to the general case  $K = 2$ , using the approach suggested by Friedman (1996).

The article is organized as follows. We introduce our model and ideas in Section 2, and in Section 3, we establish theoretical properties of our procedure. In Section 4, using a variety of classifiers, we apply our approach to simulated data and to the Lidar data mentioned above. Technical arguments are deferred to the supplementary materials.

## 2. METHODS

### 2.1 Model and Classification Problem

We observe spatial data functions  $Y_{ij}(r)$ ,  $r \in \mathcal{D}$ ,  $1 \leq i \leq n_j$ ,  $j = 1, 2$ , generated by the model

$$Y_{ij}(r) = TX_{ij}(r) + \delta_{ij}(r), \quad (2.1)$$

where  $\mathcal{D}$  denotes a  $d$ -dimensional spatial grid, or lattice;  $X_{ij}$  is the spatial function of interest,  $T$  is a linear transformation that blurs the signal; and  $\delta_{ij}$ , representing noise, is a component of a correlated stochastic process with zero mean affecting the signal. In this model, the data come from two populations,  $\Pi_1$  and  $\Pi_2$ , and, for  $j = 1, 2$ ,  $Y_{ij}$  denotes the  $i$ th data function drawn from the  $j$ th population  $\Pi_j$ , where  $i = 1, \dots, n_j$ . To simplify notation, we define scale in such a way that  $\mathcal{D} \subset \mathbb{Z}^d$ , where  $\mathbb{Z}$  is the set of all integers.

The model at Equation (2.1) is appropriate when the observations  $Y_{ij}$  are, for example, digitized images or Lidar signals. There,  $T$  typically represents the accumulated impact of issues such as generalized lens aberrations, atmospheric effects, motion blur, etc., and  $\mathcal{D}$  is a two- or three-dimensional grid.

**Remark 1:** It is important to realize that  $Y_{ij}$ ,  $X_{ij}$  and  $\delta_{ij}$  are functions defined on  $\mathcal{D}$ , and that  $T$  (and later, the transforms  $R$  and  $Q$ , which will be defined below) is not a function; it is a functional that maps the function  $X_{ij}$  to the function  $TX_{ij}$ .

Let  $Y$  be a new data value coming from  $\Pi_k$ , where  $k = 1$  or  $2$  is unknown. Our goal is to construct a classifier  $\mathcal{C}(\cdot) \equiv \mathcal{C}(\cdot | \{Y_{ij}\}_{j=1,2; i=1,\dots,n_j})$  from the data  $Y_{ij}$  that assigns  $Y$  to  $\Pi_{\hat{k}}$ , where  $\hat{k} = \mathcal{C}(Y | \{Y_{ij}\}_{j=1,2; i=1,\dots,n_j}) = 1$  or  $2$  is an estimator of  $k$ . In the applications we have in mind, where the data are images or Lidar signals, distinguishing between  $\Pi_1$  and  $\Pi_2$  is inherently a problem involving high-dimensional data analysis. In practice, the number of points  $r$  at which we observe data  $Y_{ij}(r)$  can be in the thousands, whereas the training sample sizes,  $n_1$  and  $n_2$ , are often only limited to 20 or 30.

## 2.2 Deconvolution of the Data Through the Noise Transform

As we indicated in the Introduction, when the functional  $T$  is invertible, it is sometimes argued that, instead of applying standard classifiers to the data  $Y_{ij}$ , one should apply them to inverted data where  $Y$  and each  $Y_{ij}$  are replaced by  $T^{-1}Y$  and  $T^{-1}Y_{ij}$ , or rather by regularized versions of them,  $\hat{T}^{-1}Y$  and  $\hat{T}^{-1}Y_{ij}$ . That is, classification should be based on  $\mathcal{C}(\hat{T}^{-1}Y | \{\hat{T}^{-1}Y_{ij}\}_{j=1,2; i=1,\dots,n_j})$ , instead of  $\mathcal{C}(Y | \{Y_{ij}\}_{j=1,2; i=1,\dots,n_j})$ . Transforming the data by  $T^{-1}$  is a good idea when the goal is to recover the function  $X_{ij}$ , since we have  $T^{-1}Y_{ij} = X_{ij} + T^{-1}\delta_{ij}$ , so the transformed data are no longer distorted and contain only additive noise  $T^{-1}\delta_{ij}$  of zero mean. This is only approximately true when using  $\hat{T}^{-1}$ , of course; see, for example, Cannon and Hunt (1981) and Hall (1990). However, we argue that when the goal is classification, inverting  $T$  is not necessarily a good idea, and a better strategy is to transform the data in such a way that classification performance is improved.

To explore the classification problem further, let  $\varepsilon_{ij} = TX_{ij} - T\mu_j + \delta_{ij}$ , where the function  $\mu_j$  is defined by  $\mu_j = E_f(X_{ij})$  and  $E_j$  denotes expectation, conditional on  $X_{ij}$  coming from population  $j$ . Then, (2.1) can be written as

$$Y_{ij} = T\mu_j + \varepsilon_{ij}, \quad (2.2)$$

where  $E_f(\varepsilon_{ij}) = E_f(\delta_{ij}) = 0$ . If the processes  $X_{ij} - \mu_j$  and  $\delta_{ij}$  are also linear, in particular if  $\varepsilon_{ij}$  is stationary and Gaussian, as is often approximately the case in practice, then we can write  $\varepsilon_{ij} = R\xi_{ij}$ , where  $R$  is another linear transformation and the process  $\xi_{ij}$  is white noise, that is, the random variables  $\xi_{ij}(r)$ , for  $r \in \mathbb{Z}^d$ , are uncorrelated and have zero mean and common variance  $\sigma^2$ . In this notation, the model in Equation (2.2) can be expressed as

$$Y_{ij} = T\mu_j + R\xi_{ij}, \quad (2.3)$$

so if  $R$  is invertible, Equation (2.3) can be written equivalently as

$$R^{-1}Y_{ij} = R^{-1}T\mu_j + \xi_{ij}.$$

The absence of correlation of  $\xi_{ij}$ , and the constant variances, suggests that, for a variety of classifiers, performance can be improved by working with the data  $R^{-1} Y$ , rather than with  $Y$  itself. For example, this is the case when the error process  $\varepsilon_{ij}$  in Equation (2.2) is stationary and Gaussian and we use the centroid classifier (see Section 3.1). Indeed, there, the classifier based on such transformed data is Fisher’s linear discriminant, albeit in a much higher-dimensional setting than is usually contemplated, and so, has optimality properties. In particular, this classifier is asymptotically equivalent to applying a likelihood ratio test. More generally, we shall show in Section 3 that in non-Gaussian cases, the optimal transformation, in terms of asymptotic performance of the centroid classifier, is also  $R^{-1}$ .

These considerations suggest that, for such classifiers, far from it being a good idea to replace  $Y$  and  $Y_{ij}$  by their deconvolved forms  $T^{-1} Y$  and  $T^{-1} Y_{ij}$ , we should replace them by  $R^{-1} Y$  and  $R^{-1} Y_{ij}$  and base classification on  $c(R^{-1} Y | \{R^{-1} Y_{ij}\}_{j=1,2; i=1,\dots,n_j})$ . For more general classifiers too, transforming the data prior to applying a classifier can often improve performance, but not when this transform is taken to be  $T^{-1}$ . In practice, the optimal transform is unknown and is not necessarily equal to  $R^{-1}$  for each classifier, since the best transform may depend on the particular classifier in use. Likewise, the optimal transform is not necessarily always exactly linear. However, by inverting the  $Y_{ij}$ ’s via a carefully chosen linear transform, which we shall denote by  $Q^{-1}$  in the next section, we can often improve classification performance significantly. We suggest such a practicable inversion technique in the next section, and we construct it from the data in such a way as to optimize classification performance.

### 2.3 Transforming the Data in Practice

**2.3.1 Modeling the Transform**—Since the best transform to apply to the data  $Y_{ij}$  prior to classification is generally not known, it needs to be estimated from the data. However, the sample size is usually too small for estimating this transform without imposing restrictions on it. Motivated by our discussion in the last paragraph of Section 2.2, we model the transform by the inverse  $Q^{-1}$  of a linear transform  $Q = Q_\theta$ , which depends on a low-dimensional vector of parameters  $\theta = (\theta_1, \dots, \theta_d)$ , as follows.

Let  $\omega_{Q_\theta}$  be a nonnegative weight function defined on  $\mathbb{Z}^d$  and depending on  $\theta$ . Moreover, let  $*$  denote the discrete convolution operation, defined for any two absolutely square summable functions  $f$  and  $g$  by  $f * g(r) = \sum_{s \in \mathbb{Z}^d} f(r - s) g(s)$ . We take  $Q_\theta$  to be the linear transform that maps a function  $\zeta$  to a function  $\chi_\theta = Q_\theta \zeta$  defined, for each  $r \in \mathbb{Z}^d$ , by

$$\chi_\theta(r) = \omega_{Q_\theta} * \zeta(r). \quad (2.4)$$

In image analysis terminology,  $\omega_{Q_\theta}$  is called the spread function of the transform  $Q_\theta$ . The choice of the parameters  $\theta$  will be treated in Section 2.3.3.

An example of a simple model for  $\omega_{Q_\theta}$  is the two-parameter family  $\omega_{p_0; \theta}$  where  $\theta = (\rho, \ell)$  and  $\omega_{p_0; \theta}$  is the  $\ell$ -fold convolution of the probability mass function  $p_0$ , defined by

$$p_0(r) = \left( \frac{1 - \rho}{1 + \rho} \right)^d \rho^{|r|}, \quad r \in \mathbb{Z}^d, \quad (2.5)$$

where  $|r| = \sum_{j=1}^d |r_j|$  and  $|\rho| < 1$  (usually,  $0 < \rho < 1$ ). This is the model we used in our numerical work in Section 4, but alternative models and more comments are given in Appendix A.2 in the supplementary materials.

**2.3.2 Inverting  $Q_\theta$** —Since  $Q_\theta$  is defined by a convolution, its inverse is more easily expressed in the Fourier domain. Let  $\zeta$  be a function defined on  $\mathbb{Z}^d$  such that  $\sum_{r \in \mathbb{Z}^d} |\zeta(r)| < \infty$ . The (discrete) Fourier transform  $\varphi_\zeta(t)$ , for  $t \in (-\pi, \pi)^d$ , is defined by

$$\varphi_\zeta(t) = \sum_{r \in \mathbb{Z}^d} \zeta(r) \exp(i r^T t), \quad (2.6)$$

where, on this occasion,  $i = \sqrt{-1}$ . Since the Fourier transform of a convolution between two functions is equal to the product of their Fourier transforms, we deduce from Equation (2.4) that the Fourier transform of the function  $\chi_\theta$  is given by  $\varphi_{\chi_\theta} = \varphi_\zeta \varphi_{\omega_{Q_\theta}}$

In this notation, when  $\varphi_{\omega_{Q_\theta}}(t) \neq 0$ , we can write  $\varphi_\zeta(t) = \varphi_{\chi_\theta}(t) / \varphi_{\omega_{Q_\theta}}(t)$ . If  $|\varphi_{\chi_\theta}(t) / \varphi_{\omega_{Q_\theta}}(t)|$  is integrable, then  $Q_\theta$  is invertible, and the inverse transform  $Q_\theta^{-1}$ , obtained by the Fourier inversion theorem, maps the function  $\chi_\theta$  into the function  $Q_\theta^{-1} \chi_\theta$  defined by (2.8), taking there  $\mathcal{T} = (-\pi, \pi)^d$ . If  $Q_\theta$  is not invertible, we can typically define a generalized inverse,  $Q_\theta^{-1}$ , by truncating the integral used in Fourier inversion to a small-enough set  $\mathcal{T} \subset (-\pi, \pi)^d$ , for example,

$$\mathcal{T} = \{t : \|t\| \leq \eta\} \quad \text{or} \quad \mathcal{T} = \{t : |t_j| \leq \eta, 1 \leq j \leq d\}, \quad (2.7)$$

with  $\eta \in (0, \pi)$ . Thus, in either case, we can write

$$Q_\theta^{-1} \chi_\theta(r) = (2\pi)^{-d} \int_{\mathcal{T}} \exp(-i r^T t) \{ \varphi_{\chi_\theta}(t) / \varphi_{\omega_{Q_\theta}}(t) \} dt. \quad (2.8)$$

**Remark 2:** To motivate the selections of  $\mathcal{T}$  in (2.7), observe that  $\varphi_{\omega_{Q_\theta}}(0)$  equals the sum of the weights  $\omega_{Q_\theta}(r)$  over  $r \in \mathbb{Z}^d$ , and the  $\omega_{Q_\theta}(r)$ 's would normally be chosen so that this sum was strictly positive, in fact equal to 1. Therefore,  $\varphi_{\omega_{Q_\theta}}(0) \neq 0$ , and by continuity,  $\varphi_{\omega_{Q_\theta}}(t) \neq 0$  for  $t$  in a sufficiently small neighborhood of the origin. Hence, choosing  $\mathcal{T}$  as in the formulas in (2.7), for sufficiently small  $\eta$ , ensures that the integral at (2.8) is well defined if the function  $\chi_\theta$  is uniformly bounded.

For example, if we model  $Q_\theta$  by taking  $\omega_{Q_\theta} = \omega_{p_0; \theta}$ , defined above Equation (2.5), then  $Q_\theta^{-1} \chi_\theta$  is particularly easy to calculate. As a matter of fact, by standard calculations, we have

$$\varphi_{\omega_{Q_\theta}}(t) = \varphi_{\omega_{p_0; \theta}}(t) = \prod_{j=1}^d [1 + 2\rho(1-\rho)^{-2} \{1 - \cos(t_j)\}]^{-\ell}, \quad (2.9)$$

for each  $t = (t_1, \dots, t_d)^T \in (-\pi, \pi)^d$ , so

$$Q_\theta^{-1} \chi_\theta(r) = (2\pi)^{-d} \sum_{s \in \mathbb{Z}^d} \chi_\theta(s) \int_{\mathcal{T}} \prod_{j=1}^d (e^{i(s_j - r_j)t_j} \times [1 + 2\rho(1-\rho)^{-2} \{1 - \cos(t_j)\}]^\ell) dt_j. \quad (2.10)$$

The integral in Equation (2.10) is well defined if we take  $\mathcal{T} = (-\pi, \pi)^d$ , in which case, it simplifies to

$$Q_{\theta}^{-1}\chi_{\theta}(r)=(2\pi)^{-d}\sum_{s\in\mathbb{Z}^d}\chi_{\theta}(s)\prod_{j=1}^d\int_{-\pi}^{\pi}e^{i(s_j-r_j)t_j}\times[1+2\rho(1-\rho)^{-2}\{1-\cos(t_j)\}]^{\ell}dt_j. \quad (2.11)$$

A very attractive aspect of this choice of  $Q_{\theta}$  is that we do not need smoothing parameters, such as  $\eta$  at Equation (2.7), to regularize the integral. Further, it can be proved that each integral in Equation (2.11) is equal to a constant, depending only on  $|s_j - r_j|$ ,  $\rho$ , and  $\ell$  and which vanishes if  $|s_j - r_j| > \ell$ . In other words,  $Q_{\theta}^{-1}\chi_{\theta}(r)$  is a linear combination of values of  $\chi_{\theta}(s)$ , for  $s$  in a neighborhood of  $r$  (more precisely, for  $s$  such that  $\max_{j=1,\dots,d}|s_j - r_j| \leq \ell$ ).

**2.3.3 Estimation of Unknown Parameters**—Now that we have a practicable representation  $Q_{\theta}^{-1}$  for the transform to apply to the data before classification, it remains to choose  $\theta$ . Just as, a priori, it may seem natural to invert the data by  $T^{-1}$ , it may also seem natural to choose  $\theta$  to give a good fit to the data. However, again, our goal here is to classify, and thus,  $\theta$  should rather be chosen to optimize the performance of the classifier based on  $\mathcal{C}_{\theta}(Y) \equiv \mathcal{C}(Q_{\theta}^{-1}Y|\{Q_{\theta}^{-1}Y_{ij}\}_{j=1,2;i=1,\dots,n_j})$ . We suggest choosing  $\theta$  to minimize a cross-validation estimator of error rate.

Specifically, write  $\pi_1$  for the prior probability of  $\Pi_1$ , which is typically taken to equal 0.5 if we have no a priori knowledge, or to  $n_1/(n_1 + n_2)$  if we believe that the proportion of observations from  $\Pi_1$  in the training sample is representative of that in the population. Define

$$\widehat{e}(\theta)=\frac{\pi_1}{n_1}\sum_{i=1}^{n_1}I\{\mathcal{C}_{\theta;-i1}(Y_{i1})=2\}+\frac{1-\pi_1}{n_2}\sum_{i=1}^{n_2}I\{\mathcal{C}_{\theta;-i2}(Y_{i2})=1\}, \quad (2.12)$$

where  $\mathcal{C}_{\theta;-ij}$  denotes the version of  $\mathcal{C}_{\theta}$  constructed without using  $Y_{ij}$ ; that is,  $\mathcal{C}_{\theta;-ij}(Y_{ij})=\mathcal{C}(Q_{\theta}^{-1}Y_{ij}|\{Q_{\theta}^{-1}Y_{k\ell}\}_{k=1,2;\ell=1,\dots,n_k;(k,\ell)\neq(i,j)})$ . Then,  $\widehat{e}(\theta)$  estimates the error rate,

$$e(\theta)=\pi_1 P_1\{\mathcal{C}_{\theta}(Y)=2\}+(1-\pi_1) P_2\{\mathcal{C}_{\theta}(Y)=1\}, \quad (2.13)$$

where  $P_j$  denotes probability conditional on  $Y \in \Pi_j$ . We suggest choosing  $\theta$  to minimize  $\widehat{e}(\theta)$ .

**Remark 3:** In cases where the set  $\mathcal{T}$  cannot be taken equal to  $(-\pi, \pi)^d$ , the classifier can also depend on a small number of parameters defining  $\mathcal{T}$ , which, if they are unknown, can play the role of a smoothing parameter; see the examples in (2.7). In such cases,  $\mathcal{C}_{\theta}$ ,  $\widehat{e}(\theta)$  and  $e(\theta)$  are replaced by  $\mathcal{C}_{\theta,\mathcal{T}}$ ,  $\widehat{e}(\theta, \mathcal{T})$  and  $e(\theta, \mathcal{T})$ , respectively, and  $\theta$  and  $\mathcal{T}$  re chosen to minimize  $\widehat{e}(\theta, \mathcal{T})$ .

### 3. THEORY

#### 3.1 Centroid Classifier

There exist a variety of standard classifiers that give good performance for high-dimensional data. Here, we discuss detailed theoretical properties in the context of one of the most popular and effective methods, the centroid-based technique; for example, see James and Hastie (2001) and Shin (2008). If  $\bar{Y}_j(r)=n_j^{-1}\sum_i Y_{ij}(r)$ , the centroid method assigns a new value  $Y$ , coming from  $\Pi_1$  or  $\Pi_2$ , to  $\Pi_1$  (i.e., it puts  $\mathcal{C}(Y) = 1$ ) if  $\sum_{t \in \mathcal{D}} [\{Y(t) - \bar{Y}_2(t)\}^2 -$

$(Y(r) - \bar{Y}_1(r))^2 > 0$ , and to  $\Pi_2$  (i.e., it puts  $\mathcal{C}(Y) = 2$ ) otherwise. Other classifiers will be discussed in Section 4.4.

As already highlighted in Section 2.2, if the errors are stationary, then this classifier is optimized when applied to the data  $R^{-1} Y_{ij}$ . Using the representation  $Q_\theta^{-1}$  for  $R^{-1}$ , an approximation to optimal classification involves assigning a new observation  $Y$  to  $\Pi_1$  (i.e., putting  $\mathcal{C}_\theta(Y) = 1$ ) if and only if  $S_\theta(Y) > 0$ , where

$$S_\theta(Y) = \sum_{r \in \mathcal{D}} \{|Z_\theta(r) - \bar{Z}_{2;\theta}(r)|^2 - |Z_\theta(r) - \bar{Z}_{1;\theta}(r)|^2\}, \quad (3.1)$$

with  $\bar{Z}_{j;\theta}(r) = n_j^{-1} \sum_i Z_{ij;\theta}(r)$  and where the functions  $Z_\theta$  and  $Z_{ij;\theta}$  are defined by  $Z_\theta = Q_\theta^{-1} Y$  and  $Z_{ij;\theta} = Q_\theta^{-1} Y_{ij}$ .

In this notation, the cross-validation technique for choosing  $\theta$ , described in Equation (2.12) in Section 2.3.3, can be written as

$$\widehat{c}(\theta) = \frac{\pi_1}{n_1} \sum_{i=1}^{n_1} I\{S_{\theta;-i1}(Y_{i1}) \leq 0\} + \frac{1-\pi_1}{n_2} \sum_{i=1}^{n_2} I\{S_{\theta;-i2}(Y_{i2}) > 0\}, \quad (3.2)$$

where  $S_{\theta;-ij}$  denotes the version of  $S_\theta$  in (3.1) calculated with  $\bar{Z}_j$  being replaced by  $\bar{Z}_j^{(-i)} = (n_j - 1)^{-1} \sum_{k \neq i} Z_{kj}$ . Likewise, the error rate  $e(\theta)$  in Equation (2.13) can be written as

$$e(\theta) = \pi_1 P_1\{S_\theta(Y) \leq 0\} + (1 - \pi_1) P_2\{S_\theta(Y) > 0\}. \quad (3.3)$$

### 3.2 Main Assumptions

To simplify notation, throughout Section 3, we define scale in such a way that  $\mathcal{D}$ , in  $d$ -variate Euclidean space, has edge width 1, for example,

$$\mathcal{D} = \{r = (r_1, \dots, r_d)^T : r_1, \dots, r_d \in \mathbb{Z}, -n \leq r_1, \dots, r_d \leq n\}, \quad (3.4)$$

where  $n \geq 1$ . In this setting,  $\# \mathcal{D} \asymp n^d$  and the training sample sizes,  $n_1$  and  $n_2$ , are interpreted as functions of  $n$ . Let  $\mathcal{Y}$  denote the pair of training samples  $(\mathcal{Y}_1, \mathcal{Y}_2)$ , with  $\mathcal{Y}_i = \{Y_{ij} : 1 \leq i \leq n_j\}$ ,  $Y_{ij} = (Y_{ij}(r) : r \in \mathcal{D})$ . The error rate of our classifier, computed from the training dataset  $\mathcal{Y}$ , is denoted by  $e(\theta)$  and defined in Equation (3.3). In this section, we give asymptotic formulas for  $e(\theta)$  and  $\check{e}(\theta)$ , taking  $\mathcal{T}$  to be a general subset of  $(-\pi, \pi)^d$ . For example,  $\mathcal{T}$  might be equal to  $(\pi, \pi)^d$ , or to one of the regions defined in Equation (2.7). Theory in cases where cross-validation is used to determine  $\mathcal{T}$ , as well as  $\theta$  (see Remark 3 in Section 2.3.3), can be developed at the expense of longer arguments; in the present section, we use cross-validation to optimize over  $\theta$  but not  $\mathcal{T}$ , which corresponds to our practical implementation of the method; see Section 4.

We develop our theory under three main model assumptions. First, we assume that  $R$  maps a function  $\zeta$ , defined on  $\mathbb{Z}^d$ , into a function  $R\zeta$ , defined by

$$R \zeta(r) = \omega_R * \zeta(r). \quad (3.5)$$

Second, we assume that

$$R^{-1}T\mu_j, Q_\theta^{-1}T\mu_j \text{ and } \xi \text{ are supported on } \mathcal{D}, \text{ for } j=1, 2. \quad (3.6)$$

We impose this condition only to avoid long arguments for dealing with potential edge effects. Our conclusions remain valid without it, but the proofs become considerably longer. Finally, we assume that  $T\mu_1 - T\mu_2 = T(\mu_1 - \mu_2)$  is smoother than  $\omega_R$ . More precisely, we assume that  $T(\mu_1 - \mu_2) = \alpha K * \omega_R$ , where  $\alpha$  is a constant and  $K$  is a function supported on  $\mathcal{D}$ . This assumption ensures that the inverse of the mean of the differences of the observed signals,  $R^{-1}T(\mu_1 - \mu_2)$ , remains bounded. It is imposed only to make our technical arguments simpler and explicit. If it is not satisfied, then, generally speaking, the classification problem becomes simpler, in that the difference between the means of the inverted signals is even larger and, therefore, easier to detect.

We allow the distance between the two transformed means,  $T\mu_1$  and  $T\mu_2$ , to vary with  $n$ , by letting  $\alpha$  above depend on  $n$ . In particular, we assume that  $T(\mu_1 - \mu_2) = \alpha_n K * \omega_R$ , where  $\alpha_n$  is a sequence of positive real numbers bounded above zero. The most important case is that where  $\alpha_n$  (and hence the distance) decreases with increasing  $n$ , since that enables our theoretical arguments to address particularly challenging cases. We also permit the noise variance,  $\sigma_n^2 = \text{var}\{\xi_{ij}(r)\}$ , to depend on  $n$ . We shall see that the relative sizes of  $n$ ,  $\alpha_n$ , and  $\sigma_n$  interact together to determine the performance of our classifier. Although this interaction is quite complex, to a large extent, it can be represented in terms of the quantity

$$u_n(\theta) = (\alpha_n / \sigma_n) \int_{\mathcal{D}} |\varphi_K|^2 |\varphi_{\omega_R}|^2 |\varphi_{\omega_{Q_\theta}}|^{-2} / \left\{ (2\pi)^d \int_{\mathcal{D}} |\varphi_K|^2 |\varphi_{\omega_R}|^4 |\varphi_{\omega_{Q_\theta}}|^{-4} \right\}^{1/2}, \quad (3.7)$$

where  $\varphi_{\omega_R}(t) = \sum_{r \in \mathbb{Z}^d} \omega_R(r) \exp(i r^T t)$  is the Fourier transform of  $\omega_R$  and  $\varphi_K(t) = \sum_{r \in \mathcal{D}} K(r) \exp(i r^T t)$  is the Fourier transform of  $K$ ; here, we used the fact that  $K$  is supported on  $\mathcal{D}$ .

To derive our theoretical results, we also need regularity conditions. These are more technical, and we shall describe them in detail in Appendix B.1 in the supplementary materials; see Equations (B.2)–(B.6).

### 3.3 Asymptotic Formula for Error Rate

The next theorem describes properties of  $e(\theta)$  as  $n$  diverges. Let  $\Phi$  denote the standard normal distribution function, and write  $\Theta$  for a compact set of parameters from which  $\theta$  is chosen.

**Theorem 1**—Assume that the data are generated by the model in (2.3), where  $R$  is of the form in Equation (3.5) and  $T$  is a linear transformation, and that (B.2)–(B.6) hold. Then,

$$\sup_{\theta \in \Theta} |e(\theta) - \Phi\{-u_n(\theta)\}| \rightarrow 0, \quad (3.8)$$

where the convergence is in probability.

To elucidate the implications of Theorem 1, observe first that the asymptotic error rate,  $\Phi(-u_n)$ , in Equation (3.8) is a monotone decreasing function of  $u_n$ . It therefore follows from



the formula in Equation (3.7) for  $u_n(\theta)$  that the error rate decreases as either the distance, represented by  $a_n$  between population means increases or the error variance,  $\sigma_n^2$ , decreases. Moreover, Hölder's inequality implies that  $\Phi\{-u_n(\theta)\}$ , interpreted as a functional of  $\varphi_{\omega_{Q_\theta}}$  is minimized when  $\varphi_{\omega_{Q_\theta}} = \varphi_{\omega_R}$  that is, when the transformation  $Q_\theta$  is identical to the actual transformation  $R$ .

### 3.4 Consistency of Cross-Validation Estimator of Error Rate

Recall the definition of  $\hat{\alpha}(\theta)$ , the cross-validation estimator of error rate, in Equation (3.2). Theorem 2 shows that  $\hat{\alpha}(\theta)$  shares the same asymptotic property, Equation (3.8), as the actual error rate  $\alpha(\theta)$ , and therefore, is consistent for  $\alpha(\theta)$ , uniformly in  $\theta$ .

**Theorem 2**—Assume the conditions of Theorem 1. Then,

$$\sup_{\theta \in \Theta} |\hat{\alpha}(\theta) - \alpha(\theta)| \rightarrow 0, \quad (3.9)$$

where the convergence is in probability.

Similarly, it can be proved that if  $\theta = \hat{\theta}$  is chosen to minimize  $\hat{\alpha}(\theta)$ , and used when constructing the classifier, then, under mild additional assumptions, the classifier's actual error rate will equal  $\min_{\theta \in \Theta} \alpha(\theta) + o(1)$  as  $n \rightarrow \infty$ .

## 4. NUMERICAL WORK

### 4.1 Goals of Simulations

We performed simulation studies to illustrate the following properties:

1. Transforming the data by  $T^{-1}$  prior to applying a classifier generally does not improve classification performance.
2. Transforming the data using a cross-validation-based transform  $\hat{Q}_\theta^{-1}$  generally improves classification performance even if  $Q_\theta^{-1}$  is only a rough approximation to the best transform to apply.
 

(3) The more the errors  $\epsilon_{ij}$  are correlated, the larger is the improvement at (2), especially if the error variance  $\sigma^2$  is large compared with  $T\mu_j$ .
3. The performance of classifiers, applied to data transformed by  $\hat{Q}_\theta^{-1}$ , improves as the training sample size and/or the fineness of the grid  $\mathcal{D}$  increases.

### 4.2 Simulation Setup

**4.2.1 Generation of Training Samples**—We generated training samples  $\{Y_{11}, \dots, Y_{1n_1}\}$  and  $\{Y_{21}, \dots, Y_{2n_2}\}$ , of sizes  $n_1 = n_2 = 10$  or  $n_1 = n_2 = 25$ , according to the model

$$Y_{ij}(r) = T\mu_j(r) + R\xi_{ij}(r), \quad (4.1)$$

for different curves  $\mu_j, j = 1, 2$ , and transformations  $R$  and  $T$ , and with  $r \in \mathcal{D} \subset \mathbb{R}$  or  $r = (r_1, r_2) \in \mathcal{D} \subset \mathbb{R}^2$ .

We considered four pairs of mean curves  $\mu_j$  for  $j = 1, 2$  (two univariate and two bivariate), each with several features, such as asymmetric peaks and valleys, or sinusoidal components:

- a.  $\mu_j(r) = |2r - a_j|^{4/5} \exp\{-5 \cdot 10^{-4}(4r^2 - b_j)\}$ , where  $a_1 = 5, b_1 = 100, a_2 = 4, b_2 = 80$ .

- b.  $\mu_j(r) = 9/16 \cdot c_j^{-2} (2r-50)^2 / \{1.2 + \cos(r)\}^2$ , where,  $c_1 = 200$ ,  $c_2 = 190$ .
- c.  $\mu_j(r_1, r_2) = |3r_2 - a_j|^{2/5} \exp\{-45 \cdot 10^{-4} (r_1^2 + 2r_1 r_2 + r_2^2 - b_j/9)\}$ , with  $a_j$  and  $b_j$  as in (a).
- d.  $\mu_j(r_1, r_2) = 0.1 |4 + 3r_2/50|^{1/5} \cdot \exp\{-(3r_1 + 20)/d_j\} / \{1.2 + \cos(1.5r_1)\} \cdot 1_{[-20/3, \infty)}(r_1)$ ,  
 where  $d_1 = 40$ ,  $d_2 = 50$ , and  $1_{[-20/3, \infty)}(r_1) = 1$  if  $r_1 \in [-20/3, \infty)$  and 0 otherwise.

In the previous sections, the method was discussed for a grid that had edge width 1. More generally, in our simulations, we also considered examples where the grid has edge width  $k_I$ . In that case, the various transformations have to be rescaled by a factor  $k_I$ . More precisely, if a transform  $F$  has the form  $F\xi(r) = \sum_{s \in \mathbb{Z}^d} \omega_F(s) \xi(r-s)$  on a grid of edge

width 1, on a grid of edge width  $k_I$ , it becomes  $F\xi(r) = k_I^d \sum_{s \in \mathbb{Z}_{k_I}^d} \omega_F(s) \xi(r-s)$ , where  $\mathbb{Z}_{k_I} = \{s/k_I, s \in \mathbb{Z}\}$ . Reflecting this discussion, we took  $T\mu_j(r) = k_I^d \sum_{s \in \mathbb{Z}_{k_I}^d} \omega_{\rho_0; \theta_T}(s) \mu_j(r-s)$  and

$$R\xi(r) = k_I^d \sum_{s \in \mathbb{Z}_{k_I}^d} \omega_{\rho_0; \theta_R}(s) \xi(r-s), \quad (4.2)$$

where  $\omega_{\rho_0; \theta}$  is the function defined above (2.5), with  $\theta = \theta_T(\ell_T, \rho_T)$  or  $\theta = \theta_R(\ell_R, \rho_R)$ . In our bivariate models (c) and (d), we also considered

$$R\xi(r_1, r_2) = k_I^2 \sum_{s \in \mathbb{Z}_{k_I}^2, |s_j+r_j| \leq \theta_M/k_I} \omega_M(|s_1+r_1|) \times \omega_M(|s_2+r_2|) \xi(r_1-s_1, r_2-s_2), \quad (4.3)$$

where, for  $u \in \mathbb{Z}^+$ ,  $\omega_M(u) = (\theta_M + 1 - u) / \sum_u \theta_M(\theta_M + 1 - u)$ , with  $\theta_M$  being a positive integer.

In each case, we considered several different values of  $\theta_R$  in (4.2), or  $\theta_M$  in (4.3), and we took the  $\xi_{ij}(r)$  to be independent normal  $N(0, \sigma^2)$ . Each combination of  $\sigma$  and  $\theta_R$  or  $\theta_M$  was chosen such that good classification was possible for at least one of the versions of the centroid classifier described below; see Tables A.1–A.3 in Section A.3.1, in the supplementary materials, for all the combinations we considered in practice, and for a measure of signal-to-noise ratio in each case. Finally, we took the parameter  $\theta_T$  of the transform  $T$  and the grid  $\mathcal{D}$  where the data are observed as follows:

- Model (a):  $\theta_T = (0.5, 3)$  and  $\mathcal{D} = \{-80, -80 + k_I, \dots, 80 - k_I, 80\}$ .
- Model (b):  $\theta_T = (0.5, 2)$  and  $\mathcal{D} = \{-80, -80 + k_I, \dots, 80 - k_I, 80\}$ .
- Models (c) and (d):  $\theta_T = (0.25, 2)$  and  $\mathcal{D} = \{-60, -60 + k_I, \dots, 60 - k_I, 60\} \times \{-40, -40 + k_I, \dots, 40 - k_I, 40\}$ .

In each case,  $k_I = 2$  when  $n_1 = n_2 = 10$ , and  $k_I = 1$  or 2 when  $n_1 = n_2 = 25$ . In particular, when  $n_1$  and  $n_2$  were increased, we let the grid  $\mathcal{D}$  become finer by decreasing  $k_I$  from 2 to 1 so as to illustrate point (4) in Section 4.1. We also ran simulations in the unbalanced case, where  $n_1 = 10$  and  $n_2 = 25$ , and obtained results similar to those we shall discuss below; see Figures A.4 and A.5 in Section A.3.4 in the supplementary materials.

For illustration, Figure 1 shows  $Y_{11}$  and  $Y_{12}$  in models (c) and (d), with  $\theta_R = (0.5, 3)$ . Comparing with Figure 9 in Section 4.5, we can see that model (d) looks similar to our empirical example discussed in Section 4.5.

**4.2.2 Model for  $Q_\theta$ , Generation of Test Samples and Estimation of Error Rate—**

No matter what model we used for  $R$ , we systematically modeled  $Q_\theta$  by

$$Q_\theta \xi(r) = k_{\mathcal{F}}^d \sum_{s \in \mathbb{Z}_{k,\mathcal{F}}^d} \omega_{\rho_0; \theta}(s) \xi(r-s), \quad (4.4)$$

with  $\theta = (\ell_Q, \rho_Q)$ . This model is flexible, and, as discussed in detail in Section 2.3, it has attractive practical properties such as the fact that we do not need any smoothing parameters to define  $\mathcal{T}$  in Equation (2.8), which can be taken equal to  $\mathcal{T} = (-\boldsymbol{\pi}, \boldsymbol{\pi})^d$ .

To test our classifier constructed from the training observations  $Y_{ij}$ , we generated test samples of  $N = 100$  new data curves  $Y_1^{\text{New}}, \dots, Y_{100}^{\text{New}}$ , of which half came from  $\Pi_1$  and the other half came from  $\Pi_2$ , using each time the same model as the one used to generate the  $Y_{ij}$ 's. We applied several classifiers to three versions of the  $Y_i^{\text{New}}$ 's: the untransformed noisy data  $Y_i^{\text{New}}$ , the data  $T^{-1} Y_i^{\text{New}}$ , and the data  $\hat{Q}^{-1} Y_i^{\text{New}}$ , where  $\hat{Q}$  denotes  $Q_{\hat{\theta}_{\text{CV}}}$ , with  $Q_\theta$  as in (4.4), and with  $\theta = \hat{\theta}_{\text{CV}}$  chosen to minimize the cross-validation estimator of the classification error rate, as in Section 2.3.3, where we took  $\boldsymbol{\pi}_1 = n_1 / (n_1 + n_2)$ . When  $R$  was of the form in Equation (4.2), we also applied the classifiers to the data  $R^{-1} Y_i^{\text{New}}$ .

As indicated above, we chose  $\theta = (\rho, \ell)$  to minimize the cross-validation estimator of classification error rate, where we performed the minimization over a bivariate grid of values in the range  $0 \leq \rho \leq 0.95$  and  $1 \leq \ell \leq 5$ . Here,  $\rho = 0$  denotes the identity transform, and when  $\rho = 0$ , we do not transform the data. Observe that, in our simulations and examples, the sizes of the training datasets are small, and there is little computational cost. In larger datasets, one would use  $k$ -fold cross-validation, that is, the training data would consist of a randomly selected  $(1 - k^{-1}) \times 100\%$  of the data, and the test data, the remaining  $(100/k)\%$ , with this procedure repeated many times to calculate an overall error rate. Wikipedia has a good description of this approach ([http://en.wikipedia.org/wiki/Cross-validation\\_\(statistics\)](http://en.wikipedia.org/wiki/Cross-validation_(statistics))); see also McLachlan, Do, and Ambrose (2004).

In practice, the transform  $T$  is often unknown and is not necessarily invertible. In such cases, instead of using  $T^{-1}$ , one has to use a regularized estimator  $\hat{T}^{-1}$  constructed from the data (see our real-data illustration). Here, for simplification, we take  $T$  as both known and invertible. While this may seem to be unfavorable to our approach, it actually does not matter since our point is to show that  $T^{-1}$  has essentially no role to play in our classification problem, and whether  $T^{-1}$  is known or estimated does not change our conclusions.

In each model, we generated  $B = 100$  training samples, and for each training sample, we generated a test sample of  $N = 100$  new data curves as described above, which we classified in one of the two populations using each of the methods described in the previous paragraph. For each training sample, we calculated the percentage of the new curves that were misclassified by each method. We obtained  $B = 100$  misclassification percentages for each method, and the boxplots shown in Figures 2–8 were computed from these 100 percentages.

### 4.3 Simulation Results for Centroid Classifier

**4.3.1 Data Coming From the Model in (4.1)**—We start by reporting results obtained when applying the centroid classifier described in Section 3.1 to data generated from the model in Equation (4.1). In cases where  $\hat{\epsilon}$  achieved its minimum at several values  $\theta$ , we broke the ties according to the rule described in Section A.1 in the supplementary materials. The boxplots corresponding to each of the four methods described above, for  $R$  of the form in Equation (4.2), are shown in Figures 2–4. We present the results for various values of  $\theta_R = (\rho_R, \ell_R)$ , for  $n_1 = n_2 = 10$  and  $k_I = 2$ ,  $n_1 = n_2 = 25$  and  $k_I = 2$ , and for  $n_1 = n_2 = 25$  and  $k_I = 1$ , where  $k_I$  is the distance between two adjacent univariate components of the grid  $\mathcal{D}$ . Our finite-sample results support our asymptotic theory, which implies that as  $n_1$  and  $n_2$  increase (i.e., as training sample size increases) and  $k_I$  decreases (i.e., as the grid  $\mathcal{D}$  becomes finer), the best results should be obtained by the centroid classifier applied to the data inverted by  $R^{-1}$ , of which  $Q_{\hat{\theta}_{cv}}^{-1}$  is a consistent estimator.

Overall, our results indicate that, in finite samples, it is the latter cross-validation approach that is the most competitive. This is because this method has the ability to optimize performance based on the particular sample at hand. Unsurprisingly, transforming the data through  $R^{-1}$  and  $Q_{\hat{\theta}_{cv}}^{-1}$  brings the most significant improvements when  $\rho_R$  and  $\ell_R$  are the largest, since it is in these cases that the correlation among the  $\epsilon_{ij}$ 's is the greatest. For smaller values of  $\rho_R$  and  $\ell_R$  (e.g.,  $\rho_R = 0.25$  and  $\ell_R = 1$ ), the correlation among the  $\epsilon_{ij}$ 's is relatively small, and as a result, in finite samples, the centroid method applied to the untransformed data  $Y_{ij}^{New}$  is often the most competitive approach, although even in these cases, the cross-validation approach remains highly competitive. Of course, in practice, we do not know the transformation  $R$ , and our results indicate that cross-validation-based inversion is the method of choice.

**4.3.2 Robustness Against Misspecification of R**—Next, we illustrate the robustness of the inversion procedure by reporting the results obtained when applying the centroid classifier to the data  $Q_{\hat{\theta}_{cv}}^{-1} Y_i^{New}$ , with  $Q$  as in Equation (4.4), when the true transform  $R$  was of another form, specifically the one in Equation (4.3), where we took  $\theta_M = 10, 20$ , or  $30$ . We compare this approach with the centroid-based classifier based on the data  $T^{-1} Y_i^{New}$  and with the one based on the data  $Y_i^{New}$ . We show boxplots of the percentage of misclassified data curves in Figure 5, for each of the three methods and for  $n_1 = n_2 = 10$  and  $k_I = 2$ ,  $n_1 = n_2 = 25$  and  $k_I = 2$ , and  $n_1 = n_2 = 25$  and  $k_I = 1$ . Our results indicate that even if  $Q_\theta$  at (4.4) is not the exact noise transformation, inverting the data through  $Q_{\hat{\theta}_{cv}}^{-1}$  can considerably improve on the centroid classifier based on either  $T^{-1} Y_i^{New}$  or the untransformed data  $Y_i^{New}$ .

**4.3.3 Robustness Against the Stationarity Assumption**—In practice, the model in Equation (4.1) is often an approximation to the model that generated the data. In this section, to investigate the effect of nonstationarity of the errors on our procedure, we report results of simulations where the data  $Y_{ij}$  were generated from the model

$$Y_{ij}(r) = T\mu_j(r) + R_r \xi_{ij}(r), \quad (4.5)$$

with the fixed transform  $R$  replaced by a transform  $R_r$  depending on  $r$ .

In the univariate case, instead of  $R$  in Equation (4.2), we used

$$R_r \xi(r) = k_{\mathcal{J}}^d \sum_{s \in \mathbb{Z}_{k_{\mathcal{J}}}^d} \omega_{p_0; \theta_r}(s) \xi(r-s), \quad (4.6)$$

with  $\theta_r = (\rho_r, \ell)$ , where  $\rho_r = \rho + 0.1 \cos(r/a)$  (we considered two cases:  $a = 2$  and  $a = 10$ ) and  $\rho$  and  $\ell$  are as in the previous section. In the bivariate case, instead of using the transform  $R$  in Equation (4.3) with constant  $\theta_M$ , we used the transform

$$R_r \xi(r_1, r_2) = k_{\mathcal{J}}^2 \sum_{s \in \mathbb{Z}_{k_{\mathcal{J}}}^2, |s_j+r_j| \leq \theta_{M,r_j}/k_{\mathcal{J}}} \omega_{M,r_1}(|s_1+r_1|) \omega_{M,r_2}(|s_2+r_2|) \xi(r_1-s_1, r_2-s_2), \quad (4.7)$$

where, for  $u \in \mathbb{Z}^+$  and  $j = 1, 2$ ,  $\omega_{M,r_j}(u) = (\theta_{M,r_j} + 1 - u) / \sum_u \theta_{M,r_j} (\theta_{M,r_j} + 1 - u)$ , with  $\theta_{M,r_j} = \theta_M + 2 \cdot [\alpha \cos(r_j/2)]$  (we considered two cases:  $\alpha = 2$  and  $\alpha = 4$ ),  $\theta_M$  is as in the previous section, and, for any real number  $x$ , we use  $[x]$  to denote the integer closest to  $x$ .

Although, here, the errors  $R_r \xi_{ij}(r)$  were nonstationary, we inverted the data in the same way as before, using the transform  $Q_{\hat{\theta}_{CV}}^{-1}$ . Figures 6 and 7 show boxplots of the percentage of

misclassified curves for the centroid classifier constructed from the data  $Q_{\hat{\theta}_{CV}}^{-1} Y_{ij}$ ,  $Y_{ij}$  and  $T^{-1} Y_{ij}$  where  $Y_{ij}$  was generated as in Equation (4.5), with  $\mu_j$  from model (a) and model (b), respectively,  $R_r$  as in (4.6), and  $\alpha = 2$ . For the case  $\alpha = 10$ , see Figures A.1 and A.2 in the supplementary materials. Figure 8 shows similar results for the bivariate models (c) and (d), when the data were generated according to (4.5), with  $R_r$  as in Equation (4.7) and  $\alpha = 2$ ; see Figure A.3 in the supplementary materials for the case  $\alpha = 4$ . These results indicate that our inversion method can improve classification performance significantly even when the errors are not exactly stationary; it usually does not degrade performance more than a little.

#### 4.4 Other Classifiers

Although it is beyond the scope of this article to develop theory for all types of classifiers, and derive the theoretically optimal transform for each of them, we argue that our conclusions extend to other classifiers. To illustrate this, we also implemented two other classifiers often employed in high-dimensional and functional data problems, which we

applied to the four versions of the data:  $Y_{ij}$ ,  $T^{-1} Y_{ij}$ ,  $R^{-1} Y_{ij}$  and  $Q_{\hat{\theta}_{CV}}^{-1} Y_{ij}$ , with  $\hat{\theta}_{CV}$  chosen to minimize the cross-validation estimate of classification error. Namely, we used the support vector machine (SVM) classifier with a linear kernel (svmtrain in MATLAB) and the logistic classifier applied to the partial least-square (PLS) projection of the data (here, data refer to any of the four versions, transformed or not, of the data); see Delaigle and Hall (2012b) and Section A.3.3 in the supplementary materials for more details on the logistic classifier, and see Delaigle and Hall (2012a) for properties of PLS in the functional context.

Boxplots summarizing the results of our simulations, in the same settings as for the centroid classifier, are shown in Figures A.6–A.11 in Section A.3.5 in the supplementary materials. From these figures, we can see that the results obtained with these two classifiers are very similar to those obtained with the centroid classifier. In other words, inverting by  $T^{-1}$

usually did not improve the results, and, in general, inverting by the transform  $Q_{\hat{\theta}_{CV}}^{-1}$ , chosen by cross-validation from the data, improved the results significantly (compared with using the data  $Y_{ij}$  or  $T^{-1} Y_{ij}$ ) or, when the latter worked well, transforming the data by  $Q_{\hat{\theta}_{CV}}^{-1}$  did not degrade performance much.

As already noted, the best transform to apply generally depends on the particular classifier.

However, an attractive aspect of our methodology is that the suggested inversion,  $Q_{\hat{\theta}_{CV}}^{-1}$ , is chosen to minimize a cross-validation estimator of classification error. Therefore, our approach is very flexible, since in a general setting, it approximates the inverse transform that optimizes classification.

#### 4.5 Empirical Example

We have access to data from a small experiment involving long-range infrared Lidar methods. Briefly, the idea is to discriminate between two types of aerosols that have been emitted and are to be detected by Lidar: those that are biological in nature and those that are nonbiological. There are 29 curves available to us, with  $n_1 = 15$  nonbiological and  $n_2 = 14$  biological signals.

The process involves a signal or waveform sent out in a series of bursts, and received Lidar data were observed. Some of the bursts were sent before the aerosol was released, and these were used to background-correct the received signal after the aerosol was released. For each sample, the data used here are the background-corrected received signals for a burst, 19 wavelengths, and 250 backscatter time points. In our illustrative analysis, we followed the procedure described below for 20 bursts collected almost simultaneously in the middle of the release period and then averaged over the bursts before classification. Thus, in our notation,  $Y_{ij}$  consists of the two-dimensional collection of background-corrected received signals over the wavelengths and the backscatter time points for the  $i$ th sample within the  $j$ th aerosol class. These observed data are the convolution of a true signal, the Lidar response function for a delta-pulse transmitter, with the transmitted signal. If we write  $g_{ijw}(t)$  for this true signal for wavelength  $w$  at backscatter time point  $t$ ,  $R_{ijw}(t)$  for the background-corrected received signal, and  $\mathcal{T}_w(t)$  for the transmitted signal, then, using an integral approximation to the discrete convolution, the signal we observe is

$$R_{ijw}(t) = \int_0^t \xi_{ijw}(t-v) \mathcal{T}_w(v) dv + \kappa_{ijw}(t),$$

where  $\xi_{ijw}(t)$  has mean zero. If we define  $\mathcal{N}_{jw}(t) = E\{g_{ijw}(t)\}$ ,  $\mathcal{M}_{jw}(t) = \int_0^t \mathcal{N}_{jw}(t-v) \mathcal{T}_w(v) dv$ , and  $\mathcal{P}_{ijw}(t) = \int_0^t \{\xi_{ijw}(t-v) - \mathcal{N}_{jw}(t-v)\} \mathcal{T}_w(v) dv + \kappa_{ijw}(t)$ , then we have that the observed data are given by  $R_{ijw}(t) = \mathcal{M}_{jw}(t) + \mathcal{P}_{ijw}(t)$ , where  $\mathcal{P}_{ijw}(t)$  has mean zero. In our notation,  $Y_{ij}$ ,  $\mu_j$ ,  $T$ ,  $\mu_j$ , and  $\epsilon_{ij}$  are the collection of  $R_{ijw}(t)$ ,  $\mathcal{N}_{jw}(t)$ ,  $\mathcal{M}_{jw}(t)$ , and  $\mathcal{P}_{ijw}(t)$  over the wavelengths and backscatter ranges, respectively, but averaged across 20 bursts. It is readily observed that the transformation  $T\mu_j$  is linear. Two typical observed average curves for each population are given in Figure 9.

We considered three approaches. The first simply used the observed data  $Y_{ij}$ . The second was our method applied to the  $Q_{\hat{\theta}_{CV}}^{-1} Y_{ij}$ , where  $Q_{\theta}$  had the form as in Equation (4.4). In the third, for each burst and each wavelength, we deconvolved to estimate  $g_{ijw}(t)$  using the Wiener–Helstrom method described by Warren et al. (2008), and averaged over the bursts. In each case, since we could not generate new data, we estimated the misclassification error rate (i.e., misclassification percentage) by cross-validation. In other words, as in the case of the procedure described in Section 2.3.3, we built the classifier from all but one of the 29 curves, classified that curve in one of the two populations (nonbiological or biological), and averaged the results over all 29 curves.

For the centroid classifier, the cross-validation estimator of the misclassification error rate was 34.5% for the first approach based on nontransformed data, 24.1% for our cross-validation-based inversion approach, and 34.5% for the third approach based on inversion of  $T$ . For the SVM and logistic regression classifiers, the estimator of the misclassification error rate was 37.9% (SVM) or 27.6% (logistic) when the classifier was based on nontransformed data, 17.2% (SVM) or 21% (logistic) when the classifier was based on our cross-validation-based inversion method, and 58.6% (SVM) or 31% (logistic) when the classifier was based on inversion of  $T$ . For all three classifiers, the reduction in the misclassification error rate obtained by our cross-validation-based data inversion illustrates the significant improvement that can be obtained by inverting the data through a data-driven transform chosen to minimize an estimator of classification error.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

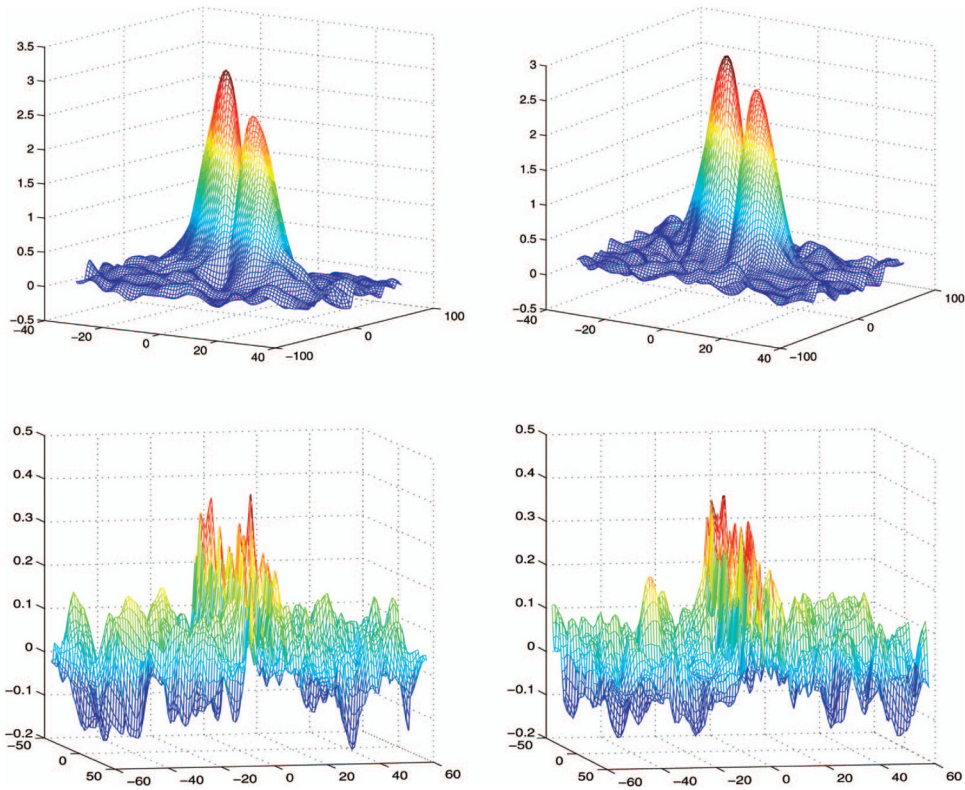
Carroll's research was supported by a grant from the National Cancer Institute (R37-CA057030) and in part by award number KUS-CI-016-04, made by King Abdullah University of Science and Technology (KAUST) and by the National Science Foundation (DMS-0914951). Delaigle's research was supported by grants and a Queen Elizabeth II Fellowship from the Australian Research Council, and Hall's research was supported by a Federation Fellowship, a Laureate Fellowship, and grants from the Australian Research Council.

## References

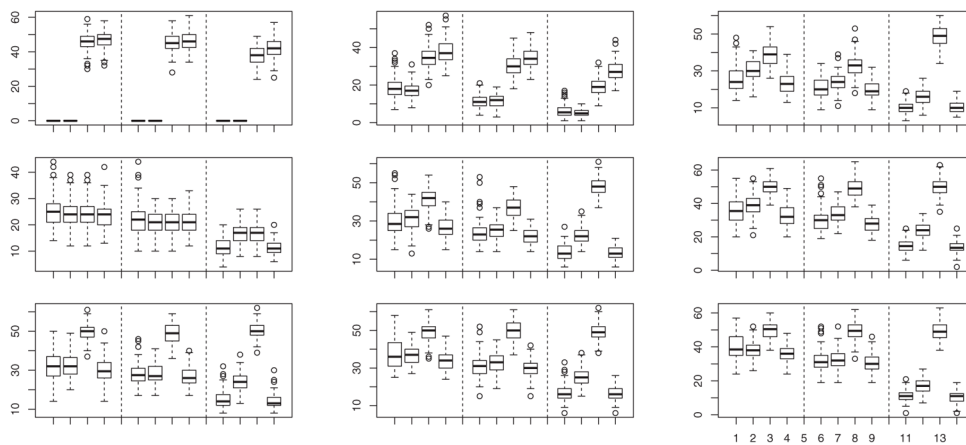
- Besag J. On the Statistical-Analysis of Dirty Pictures. *Journal of the Royal Statistical Society, Series B*. 1986; 48:259–302.
- Cannon M. Blind Deconvolution of Spatially Invariant Image Blurs With Phase. *IEEE Transactions on Acoustics, Speech and Signal Processing*. 1976; 24:58–63.
- Cannon TM, Hunt BR. *Image Processing by Computer*. Scientific American. 1981; 245:214–225.
- Carasso AS. Direct Blind Deconvolution. *SIAM Journal on Applied Mathematics*. 2001; 61:1980–2007.
- Cressie N, Kornak J. Spatial Statistics in the Presence of Location Error With an Application to Remote Sensing of the Environment. *Statistical Science*. 2003; 18:436–456.
- Crosilla F, Visintini D, Sepic F. An Automatic Classification and Robust Segmentation Procedure of Spatial Objects. *Statistical Methods and Applications*. 2007; 15:329–341.
- Dass SC, Nair VN. Edge Detection, Spatial Smoothing, and Image Reconstruction with Partially Observed Multivariate Data. *Journal of the American Statistical Association*. 2003; 98:77–89.
- Delaigle A, Hall P. Methodology and Theory for Partial Least Squares Applied to Functional Data. *The Annals of Statistics*. 2012a; 40:322–352.
- Delaigle A, Hall P. Componentwise Classification and Clustering of Functional Data. *Biometrika*. 2012b; 99(2):299–313.
- Donoho DL. Statistical Estimation and Optimal Recovery. *The Annals of Statistics*. 1994; 22:238–270.
- Figueiredo MAT, Nowak RD. An EM Algorithm for Wavelet-Based Image Restoration. *IEEE Transactions on Image Processing*. 2003; 12:906–916. [PubMed: 18237964]
- Friedman, JH. Another Approach to Polychotomous Classification. Technical Report. 1996. available at <http://www-stat.stanford.edu/jhf/ftp/poly.pdf>
- Galatsanos NP, Mesarovi VZ, Molina R, Katsaggelos AK, Mateos J. Hyperparameter Estimation in Image Restoration Problems With Partially Known Blurs. *Optical Engineering*. 2002; 41:1845–1854.
- Hall P. Optimal Convergence Rates in Signal Recovery. *The Annals of Probability*. 1990; 18:887–900.
- Hall P, Qiu P. Blind Deconvolution and Deblurring in Image Analysis. *Statistica Sinica*. 2007a; 17:1483–1509.

- Hall P, Qiu P. Nonparametric Estimation of a Point Spread Function in Multivariate Problems. *The Annals of Statistics*. 2007b; 35:1512–1534.
- Huang HC, Cressie N. Deterministic/Stochastic Wavelet Decomposition for Recovery of Signal From Noisy Data. *Technometrics*. 2000; 42:262–276.
- Huang XF, Qiu P. Blind Deconvolution for Jump-Preserving Curve Estimation. *Mathematical Problems in Engineering* [online]. 2010:Article No. 350849.10.1155/2010/350849
- James G, Hastie T. Functional Linear Discriminant Analysis for Irregularly Sampled Curves. *Journal of the Royal Statistical Society, Series B*. 2001; 63:533–550.
- Johnstone IM. Speed of Estimation in Positron Emission Tomography and Related Inverse Problems. *The Annals of Statistics*. 1990; 18:251–280.
- Joshi MV, Chaudhuri S. Joint Blind Restoration and Surface Recovery in Photometric Stereo. *Journal of the Optical Society of America*. 2005; A22:1066–1076. [PubMed: 15984479]
- Klein R, Press SJ. Adaptive Bayesian Classification of Spatial Data. *Journal of the American Statistical Association*. 1992; 87:844–851.
- Kundur D, Hatzinakos D. A Novel Blind Deconvolution Scheme for Image Restoration Using Recursive Filtering. *IEEE Transactions on Signal Processing*. 1998; 46:375–389.
- McLachlan, GJ.; Do, KA.; Ambrose, C. *Analyzing Microarray Gene Expression Data*. Hoboken, NJ: Wiley; 2004.
- Mukherjee PS, Qiu P. 3-D Image Denoising by Local Smoothing and Nonparametric Regression. *Technometrics*. 2011; 53:196–208.
- Popescu DC, Hellicar AD. Point Spread Function Estimation for a Terahertz Imaging System. *EURASIP Journal on Advances in Signal Processing* [online]. 2010:Article No. 575817.10.1155/2010/575817
- Qiu, P. *Image Processing and Jump Regression Analysis*. New York: Wiley; 2005.
- Qiu P. Jump Surface Estimation, Edge Detection, and Image Restoration. *Journal of the American Statistical Association*. 2007; 102:745–756.
- Qiu P. A Nonparametric Procedure for Blind Image Deblurring. *Computational Statistics and Data Analysis*. 2008; 52:4828–4841.
- Shi T, Cressie N. Global Statistical Analysis of MISR Aerosol Data: A Massive Data Product From NASA's Terra Satellite. *Environ-metrics*. 2007; 18:665–680.
- Shin H. An Extension of Fisher's Discriminant Analysis for Stochastic Processes. *Journal of Multivariate Analysis*. 2008; 99:1191–1216.
- Warren RE, Vanderbeek RG, Ben-David A, Ahl JL. Simultaneous Estimation of Aerosol Cloud Concentration and Spectral Backscatter From Multiple-Wavelength Lidar Data. *Applied Optics*. 2008; 47:4309–4320. [PubMed: 18716635]

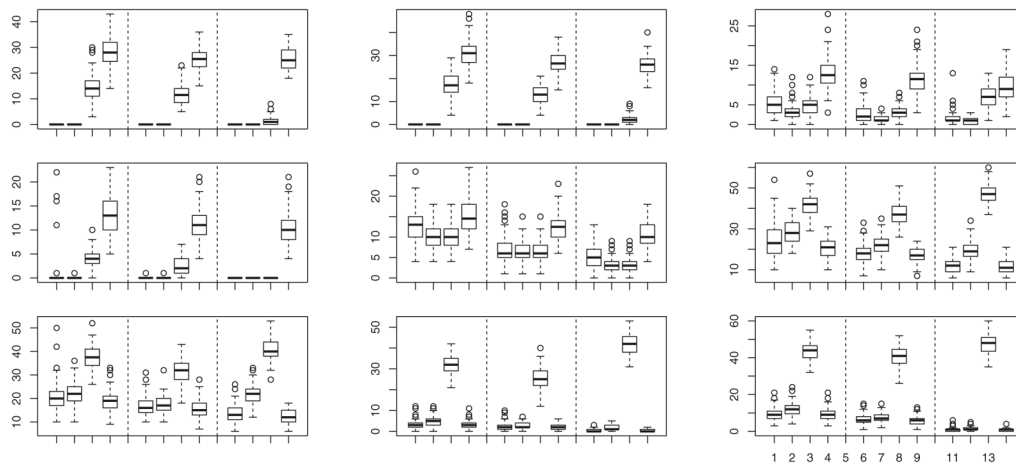




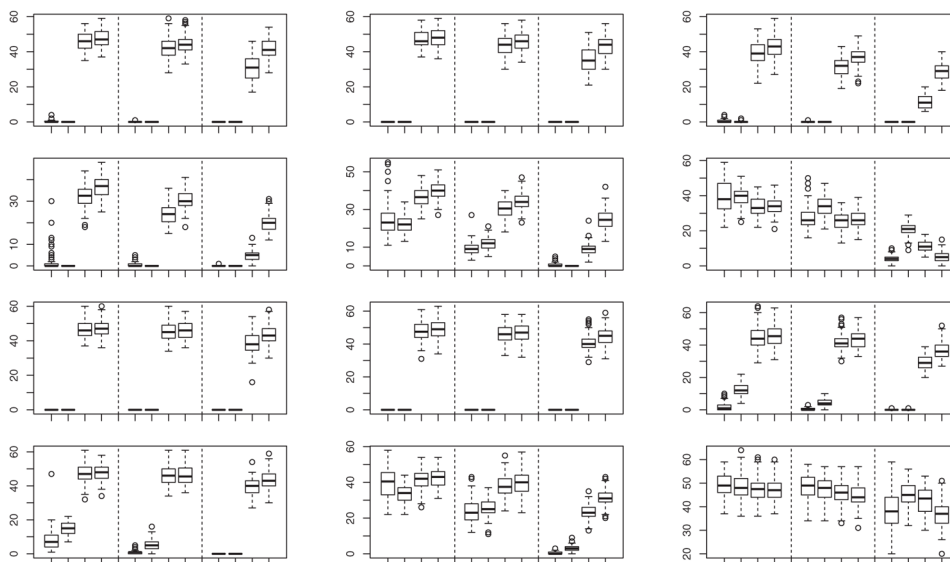
**Figure 1.** Plots of  $Y_{11}$  (left) and  $Y_{12}$  (right), in models (c) (row 1) and (d) (row 2) with  $\theta_R = (0.5, 3)$ . The online version of this figure is in color.



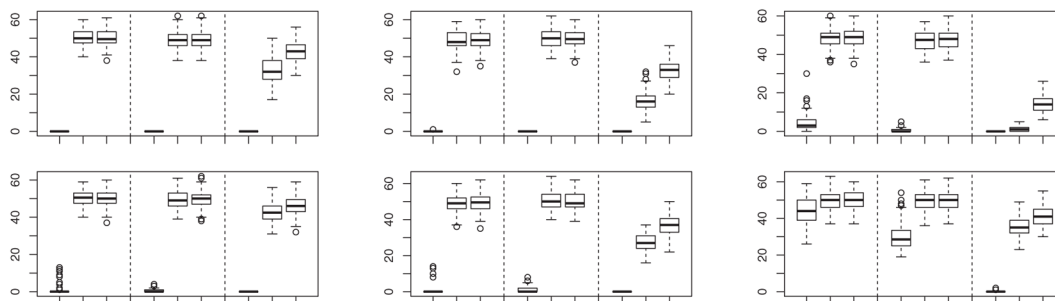
**Figure 2.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from model (a) when  $\theta_R = (\rho_R, \ell_R)$ , with  $\rho_R = 0.75, 0.5,$  and  $0.25$  in rows 1, 2, and 3, respectively; and  $\ell_R = 3, 2,$  and  $1$  in columns 1, 2, and 3, respectively. In each group of 12 boxes, the first four are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next four are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last four are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of four boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box),  $R^{-1}$  (second box), or  $T^{-1}$  (third box), or are untransformed (fourth box).



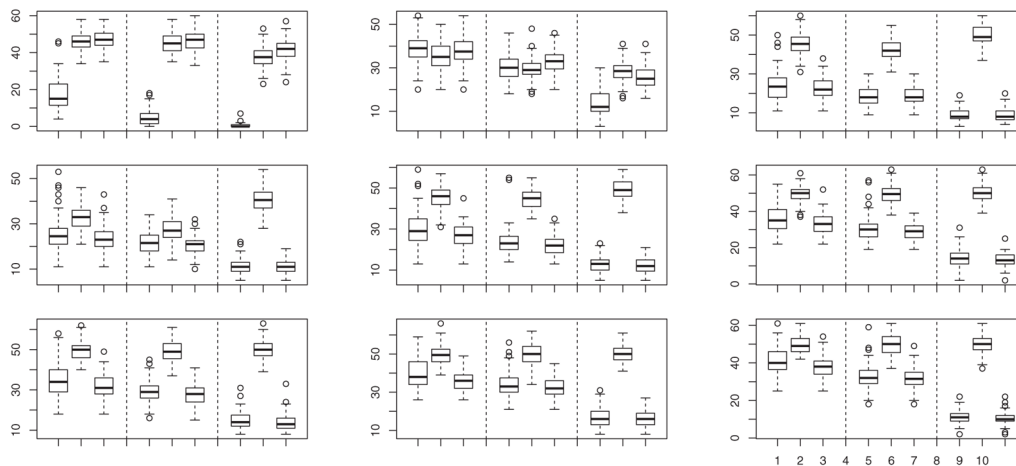
**Figure 3.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from model (b) when  $\theta_R = (\rho_R, \ell_R)$ , with  $\rho_R = 0.75, 0.5,$  and  $0.25$  in rows 1, 2, and 3, respectively; and  $\ell_R = 3, 2,$  and  $1$  in columns 1, 2, and 3, respectively. In each group of 12 boxes, the first four are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next four are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last four are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of four boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box),  $R^{-1}$  (second box), or  $T^{-1}$  (third box), or are untransformed (fourth box).



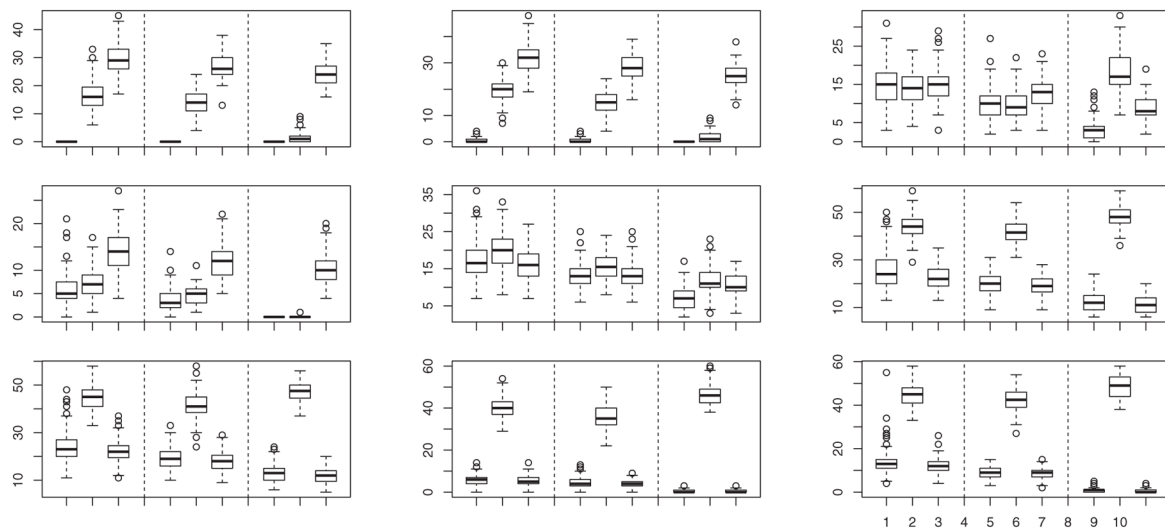
**Figure 4.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from models (c) (rows 1 and 2) and (d) (rows 3 and 4) when  $\theta_R = (\rho_R, \ell_R)$ , with  $\rho_R = 0.85$  and  $0.5$  in rows 1, 3 and 2, 4, respectively; and  $\ell_R = 3, 2,$  and  $1$  in columns 1, 2, and 3, respectively. In each group of 12 boxes, the first four are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next four are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last four are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of four boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box),  $R^{-1}$  (second box), or  $T^{-1}$  (third box), or are untransformed (fourth box).



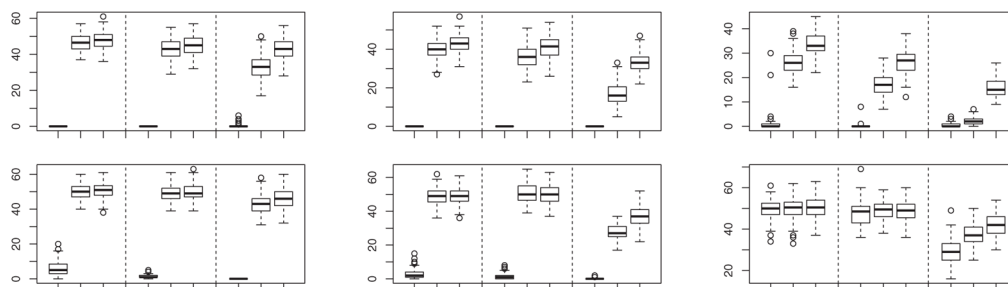
**Figure 5.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from models (c) (row 1) and (d) (rows 2) when  $R$  is of the form at (4.3), with  $\theta_M = 30, 20,$  and  $10$  in columns 1, 2, and 3, respectively. In each group of nine boxes, the first three are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next three are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last three are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of three boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box) or  $T^{-1}$  (second box), or are untransformed (third box).



**Figure 6.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from model (a) when  $\theta_R = (\rho_{R,I}, l_R)$ , with  $\rho_{R,I} = \rho_R + 0.1 \cos(\pi/2)$ ;  $\rho_R = 0.75, 0.5,$  and  $0.25$  in rows 1, 2, and 3, respectively; and  $l_R = 3, 2,$  and  $1$  in columns 1, 2, and 3, respectively. In each group of nine boxes, the first three are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next three are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last three are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of three boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box) or  $T^{-1}$  (second box), or are untransformed (third box).

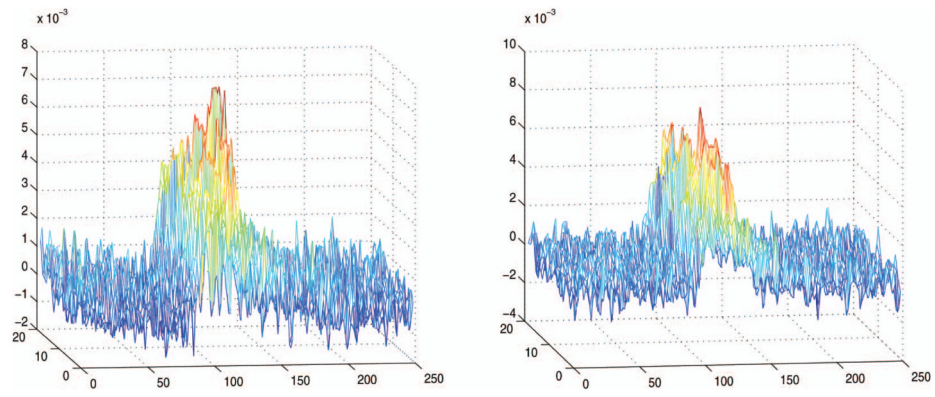


**Figure 7.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from model (b) when  $\theta_R = (\rho_{R,i}, \ell_R)$ , with  $\rho_{R,i} = \rho_R + 0.1 \cos(i\pi/2)$ ;  $\rho_R = 0.75, 0.5$ , and  $0.25$  in rows 1, 2, and 3, respectively; and  $\ell_R = 3, 2$ , and  $1$  in columns 1, 2, and 3, respectively. In each group of nine boxes, the first three are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next three are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last three are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of three boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box) or  $T^{-1}$  (second box), or are untransformed (third box).



**Figure 8.** Boxplots of percentage of misclassified observations calculated from 100 simulated samples from models (c) (row 1) and (d) (rows 2) when  $R_r$  is of the form at (4.7), with  $\theta_{M,rj} = \theta_M + 2 \cdot [2 \cos(rj/2)]$  and  $\theta_M = 30, 20,$  and  $10$  in columns 1, 2, and 3, respectively. In each group of nine boxes, the first three are for  $n_1 = n_2 = 10$  and  $k_I = 2$ , the next three are for  $n_1 = n_2 = 25$  and  $k_I = 2$ , and the last three are for  $n_1 = n_2 = 25$  and  $k_I = 1$ . In each group of three boxes, the data are transformed by  $\hat{Q}^{-1}$  (first box) or  $T^{-1}$  (second box), or are untransformed (third box).





**Figure 9.** Plots of two background-corrected received data curves from population 1 (left) and population 2 (right), averaged over the 20 bursts, across wavelength and the backscatter spectral range. The online version of this figure is in color.