

## TWO-SAMPLE TESTS IN FUNCTIONAL DATA ANALYSIS STARTING FROM DISCRETE DATA

Peter Hall<sup>1,2</sup> and Ingrid Van Keilegom<sup>2</sup>

<sup>1</sup>*The University of Melbourne* and <sup>2</sup>*Université catholique de Louvain*

*Abstract:* One of the ways in which functional data analysis differs from other areas of statistics is in the extent to which data are pre-processed prior to analysis. Functional data are invariably recorded discretely, although they are generally substantially smoothed as a prelude even to viewing by statisticians, let alone further analysis. This has a potential to interfere with the performance of two-sample statistical tests, since the use of different tuning parameters for the smoothing step, or different observation times or subsample sizes (i.e., numbers of observations per curve), can mask the differences between distributions that a test is trying to locate. In this paper, and in the context of two-sample tests, we take up this issue. Ways of pre-processing the data, so as to minimise the effects of smoothing, are suggested. We show theoretically and numerically that, by employing exactly the same tuning parameter (e.g. bandwidth) to produce each curve from its raw data, significant loss of power can be avoided. Provided a common tuning parameter is used, it is often satisfactory to choose that parameter along conventional lines, as though the target was estimation of the continuous functions themselves, rather than testing hypotheses about them. Moreover, in this case, using a second-order smoother (such as a local-linear method), the subsample sizes can be almost as small as the square root of sample sizes before the effects of smoothing have any first-order impact on the results of a two-sample test.

*Key words and phrases:* Bandwidth, bootstrap, curve estimation, hypothesis testing, kernel, Cramér-von Mises test, local-linear methods, local-polynomial methods, nonparametric regression, smoothing.

### 1. Introduction

Although, in functional data analysis, the data are treated as though they are in the form of curves, in practice they are invariably recorded discretely. They are subject to a pre-processing step, usually based on local-polynomial or spline methods, to transform them to the smooth curves to which the familiar algorithms of functional data analysis are applied. In many instances the pre-processing step is not of great importance. However, in the context of two-sample hypothesis testing it has the potential to significantly reduce power. Our aim in the present paper is to explore this issue, and suggest methods which allow the effects of smoothing to be minimised.

Although smoothing methods are sometimes used in related contexts, for example to produce alternatives to traditional two-sample hypothesis tests, the mathematical models with which a statistician typically works allow for the discrete, unprocessed form of the data. See, for example, Dette and Neumeyer (2001). Moreover, in such instances the data are passed from one statistician to another before processing, i.e., before smoothing.

By way of contrast, functional datasets are usually exchanged by researchers in post-processed form, after the application of a smoothing algorithm; that is seldom the case for data in related problems such as nonparametric regression. The widespread use of pre-process smoothing for functional data makes the effects of smoothing more insidious than usual, and strengthens motivation for understanding the impact that smoothing might have.

There is rightly a debate as to whether statistical smoothing should be used at all, in a conventional sense, when constructing two-sample hypothesis tests. Employing a small bandwidth (in theoretical terms, a bandwidth which converges to zero as sample size increases) can reduce statistical power unnecessarily, although from other viewpoints power can be increased. See, for example, the discussion by Ingster (1993), Fan (1994), Anderson, Hall and Titterington (1994), Fan (1998), Fan and Ullah (1999) and Li (1999).

It has also been observed, in a range of settings, that when probability statements rather than  $L_p$  metrics are to be optimised, undersmoothing can give improved performance. In particular, the level accuracy of tests can be improved by undersmoothing; see, for example, Hall and Sheather (1988), Koenker and Machado (1999) and Koenker and Xiao (2002).

In the context of two-sample hypothesis testing, our main recommendation is appropriate in cases where the “subsample sizes” (that is, the numbers of points at which data are recorded for individual functions) do not vary widely. There we suggest that exactly the same tuning parameters be used to produce each curve from its raw data, for all subsamples in both datasets. For example, when using kernel-based methods this would mean using the same bandwidth in all cases; for splines it would mean using the same penalty, or the same knots. Such a choice ensures that, under the null hypothesis that the two samples of curves come from identical populations, the main effects of differing observation-time distributions and differing subsample sizes (for the different curves) cancel. As a result, the effect that smoothing has on bias is an order of magnitude less than it would be if different bandwidths, tuned to the respective curves, were used. The latter approach can lead to power loss for a two-sample test.

The fact that the common tuning parameter can be chosen by a conventional curve-estimation method, such as cross-validation or a plug-in rule, is an attractive feature of the proposal. New smoothing-parameter choice algorithms are not

essential. Nevertheless, it can be shown that under more stringent assumptions a bandwidth of smaller size can be advantageous. See Sections 3.3 and 4, and also Hall and Van Keilegom (2006), for discussion.

In a more general setting, where  $k$  samples are available and our task is to test the null hypothesis that all are identical, a test is typically constructed by taking either the weighted sum, or the maximum, of values of a statistic  $\widehat{T}_{j_1 j_2}$  that tests for differences between populations  $j_1$  and  $j_2$ . However, if the test leads to rejection, then one generally seeks further information about differences among populations that pairwise tests can provide. For these reasons, the conclusions we reach for two-sample hypothesis tests have direct implications for  $k$ -sample problems.

Recent work on two-sample hypothesis tests in nonparametric and semiparametric settings includes that of Louani (2000), Claeskens, Jing, Peng and Zhou (2003) and Cao and Van Keilegom (2006). Extensive discussion of methods and theory for functional data analysis is given by Ramsay and Silverman (1997, 2002). Recent contributions to hypothesis testing in this field include those of Fan and Lin (1998), Locantore, Marron, Simpson, Tripoli, Zhang and Cohn (1999), Spitzner, Marron and Essick (2003), Cuevas, Febrero and Fraiman (2004) and Shen and Faraway (2004). However, this work does not broach the subject of the effect that pre-processing has on results.

## 2. Statement of Problem and Methodology

### 2.1. The data and the problem.

We observe data

$$\begin{aligned} U_{ij} &= X_i(S_{ij}) + \delta_{ij}, & 1 \leq i \leq m, & \quad 1 \leq j \leq m_i, \\ V_{ij} &= Y_i(T_{ij}) + \epsilon_{ij}, & 1 \leq i \leq n, & \quad 1 \leq j \leq n_i, \end{aligned} \quad (2.1)$$

where  $X_1, X_2, \dots$  are identically distributed as  $X$ ;  $Y_1, Y_2, \dots$  are identically distributed as  $Y$ ; the  $\delta_{ij}$ 's are identically distributed as  $\delta$ ; the  $\epsilon_{ij}$ 's are identically distributed as  $\epsilon$ ;  $X$  and  $Y$  are both random functions, defined on the interval  $\mathcal{I} = [0, 1]$ ; the observation errors,  $\delta_{ij}$  and  $\epsilon_{ij}$ , have zero means and uniformly bounded variances; the sequences of observation times,  $S_{i1}, \dots, S_{im_i}$  and  $T_{i1}, \dots, T_{in_i}$ , are either regularly spaced on  $\mathcal{I}$  or drawn randomly from a distribution (possibly different for each  $i$ , and also for  $S$  and  $T$ ) having a density that is bounded away from zero on  $\mathcal{I}$ ; and the quantities  $X_{i_1}, Y_{i_2}, S_{i_1 j_1}, T_{i_2 j_2}, \delta_{i_1}$  and  $\epsilon_{i_2}$ , for  $1 \leq i_1 \leq m, 1 \leq j_1 \leq m_{i_1}, 1 \leq i_2 \leq n$  and  $1 \leq j_2 \leq n_{i_2}$ , are all totally independent.

Given the data at (2.1), we wish to test the null hypothesis,  $H_0$ , that the distributions of  $X$  and  $Y$  are identical. In many cases of practical interest,  $X$  and

$Y$  would be continuous with probability 1, and then  $H_0$  would be characterised by the statement,

$$F_X(z) = F_Y(z) \quad \text{for all continuous functions } z,$$

where  $F_X$  and  $F_Y$  are the distributional functionals of  $X$  and  $Y$ , respectively:

$$\begin{aligned} F_X(z) &= P\{X(t) \leq z(t) \quad \text{for all } t \in \mathcal{I}\}, \\ F_Y(z) &= P\{Y(t) \leq z(t) \quad \text{for all } t \in \mathcal{I}\}. \end{aligned}$$

Thus, our approach to analysis will focus on differences between distributions, rather than (for instance) on differences between means, of random functions. However, conclusions similar to those discussed here would be drawn in the context of testing hypotheses about a mean. Issues connected with equality of distributions are increasingly of interest even in parametric settings, for example in connection to the fitting of models (e.g., the equilibrium-search model of Mortensen (1990) and Burdett and Mortensen (1998)) to wage and employment data. In such problems, goodness of fit can be addressed by drawing a very large (effectively, infinite) sample from the population defined by the model with parameters fitted, and applying a two-sample test in the context of those data and the sample that is being compared with the model.

In some problems, the time points  $S_{ij}$  and  $T_{ij}$  would be regularly spaced on grids of the same size. For example, they might represent the times of monthly observations of a process, such as the number of tons of a certain commodity exported by a particular country in month  $j$  of year  $i$ , in which case  $m_i = n_i = 12$  for each  $i$ . (The differing lengths of different months could usually be ignored. In examples such as this we might wish first to correct the data for linear or periodic trends.) Here the assumption of independence might not be strictly appropriate, but the methods that we suggest will be approximately correct under conditions of weak dependence.

In other cases the values of  $m_i$  and  $n_i$  can vary from one index  $i$  to another, and in fact those quantities might be modelled as conditioned values of random variables. The observation times may also exhibit erratic variation. For example, in longitudinal data analysis,  $U_{ij}$  might represent a measurement of the condition of the  $i$ th type- $X$  patient at the  $j$ th time in that patient's history. Since patients are seen only at times that suit them, then both the values of the observation times, and the number of those times, can vary significantly from patient to patient.

In general, unless we have additional knowledge about the distributions of  $X$  and  $Y$  (for example, both distributions are completely determined by finite parameter vectors), we cannot develop a theoretically consistent test unless the

values of  $m_i$  and  $n_i$  increase without bound as  $m$  and  $n$  increase. Therefore, the divergence of  $m_i$  and  $n_i$  will be a key assumption in our theoretical analysis.

**2.2. Methodology for testing hypotheses**

Our approach to testing  $H_0$  is to compute estimators,  $\widehat{X}_i$  and  $\widehat{Y}_i$ , of  $X_i$  and  $Y_i$ , by treating the problem as one of nonparametric regression and passing nonparametric smoothers through the datasets  $(S_{i1}, U_{i1}), \dots, (S_{im_i}, U_{im_i})$  and  $(T_{i1}, V_{i1}), \dots, (T_{in_i}, V_{in_i})$ , respectively. Then, treating the functions  $\widehat{X}_1, \dots, \widehat{X}_m$  and  $\widehat{Y}_1, \dots, \widehat{Y}_n$  as independent and identically distributed observations of  $X$  and  $Y$ , respectively (under our assumptions they are at least independent), we construct a test of  $H_0$ .

For example, we might compute estimators  $\widehat{F}_X$  and  $\widehat{F}_Y$  of  $F_X$  and  $F_Y$ , respectively:

$$\widehat{F}_X(z) = \frac{1}{m} \sum_{i=1}^m I(\widehat{X}_i \leq z), \quad \widehat{F}_Y(z) = \frac{1}{n} \sum_{i=1}^n I(\widehat{Y}_i \leq z), \quad (2.2)$$

where the indicator  $I(\widehat{X}_i \leq z)$  is interpreted as  $I\{\widehat{X}_i(t) \leq z(t) \text{ for all } t \in \mathcal{T}\}$ , and  $I(\widehat{Y}_i \leq z)$  is interpreted similarly. These quantities might be combined into a test statistic of Cramér-von Mises type, say

$$\widehat{T} = \int \{\widehat{F}_X(z) - \widehat{F}_Y(z)\}^2 \mu(dz), \quad (2.3)$$

where  $\mu$  denotes a probability measure on the space of continuous functions.

An alternative approach would be to use, rather than  $\widehat{T}$ , a Kolmogorov-Smirnov statistic,  $\widehat{T}' = \sup_z |\widehat{F}_X(z) - \widehat{F}_Y(z)|$ . However, this would typically produce a test with less power than a test founded on  $\widehat{T}$ , since for the latter statistic the differences between the two distribution functions are averaged over all  $z$  values for which they arise, rather than just in the vicinity of the value of  $z$  for which the distance is greatest.

The integral in (2.3) can be calculated by Monte Carlo simulation, for example as

$$\widehat{T}_N = \frac{1}{N} \sum_{i=1}^N \{\widehat{F}_X(M_i) - \widehat{F}_Y(M_i)\}^2, \quad (2.4)$$

where  $M_1, \dots, M_N$  are independent random functions with the distribution of  $M$ , say, for which  $\mu(A) = P(M \in A)$  for each Borel set  $A$  in the space of continuous functions on  $\mathcal{T}$ . Of course, the  $M_i$ 's are independent of  $\widehat{F}_X$  and  $\widehat{F}_Y$ , and  $\widehat{T}_N \rightarrow \widehat{T}$ , with probability one conditional on the data, as  $N \rightarrow \infty$ . Note that in the latter result,  $\widehat{T}$  does not depend on  $N$ .

### 2.3. Methodology for estimating $X_i$ and $Y_i$

For brevity we confine attention to just one technique, local-polynomial methods, for computing  $\widehat{X}_i$  and  $\widehat{Y}_i$ . (Results can be expected to be similar if one uses other conventional smoothers, for example splines.) Taking the degree of the polynomial to be odd, and estimating  $X_i$ , we compute the value  $(\widehat{a}_0, \dots, \widehat{a}_{2r+1})$  of the vector  $(a_0, \dots, a_{2r+1})$  that minimises

$$\sum_{j=1}^{m_i} \left\{ U_{ij} - \sum_{k=0}^{2r-1} a_k (S_{ij} - t)^k \right\}^2 K \left( \frac{t - S_{ij}}{h_{X_i}} \right),$$

where  $r \geq 1$  is an integer,  $h_{X_i}$  is a bandwidth, and  $K$ , the kernel function, is a bounded, symmetric, compactly supported probability density. Then,  $\widehat{a}_0 = \widehat{X}_i(t)$ .

In the particular case  $r = 1$  we obtain a local-linear estimator of  $X_i(t)$ ,

$$\widehat{X}_i(t) = \frac{A_{i2}(t) B_{i0}(t) - A_{i1}(t) B_{i1}(t)}{A_{i0}(t) A_{i2}(t) - A_{i1}(t)^2},$$

where

$$A_{ir}(t) = \frac{1}{m_i h_{X_i}} \sum_{j=1}^{m_i} \left( \frac{t - S_{ij}}{h_{X_i}} \right)^r K \left( \frac{t - S_{ij}}{h_{X_i}} \right),$$

$$B_{ir}(t) = \frac{1}{m_i h_{X_i}} \sum_{j=1}^{m_i} U_{ij} \left( \frac{t - S_{ij}}{h_{X_i}} \right)^r K \left( \frac{t - S_{ij}}{h_{X_i}} \right),$$

$h_{X_i}$  denotes a bandwidth and  $K$  is a kernel function. The estimator  $\widehat{Y}_i$  is constructed similarly. Local-linear methods have an advantage over higher-degree local-polynomial approaches in that they suffer significantly less from difficulties arising from singularity, or near-singularity, of estimators.

Treating  $X_i$  as a fixed function (that is, fixing  $i$  and conditioning on the stochastic process  $X_i$ ), assuming that  $X_i$  has  $2(r+1)$  bounded derivatives, and that  $h_{X_i}$  is chosen of size  $m_i^{-1/(2r+1)}$ , the estimator  $\widehat{X}_i$  converges to  $X_i$  at the mean-square optimal rate  $m_i^{-2r/(2r+3)}$ , as  $m_i$  increases. See, for example, Fan (1993), Fan and Gijbels (1996) and Ruppert and Wand (1994) for discussion of both practical implementation and theoretical issues.

### 2.4. Bandwidth choice

A number of potential bandwidth selectors are appropriate when all the subsample sizes  $m_i$  and  $n_j$  are similar and the bandwidths  $h = h_{X_i} = h_{Y_j}$  are identical. Theoretical justification for using a common bandwidth, when the goal is hypothesis testing rather than function estimation, will be given in Section 3.

One approach to common bandwidth choice is to use a “favourite” method to compute an empirical bandwidth for each curve  $X_i$  and  $Y_j$ , and then take the average value to be the common bandwidth. Another technique, appropriate in the case of plug-in rules, is to use an average value of each of the components of a plug-in bandwidth selector, and assemble the average values, using the plug-in formula, to form the common bandwidth. A third approach, valid in the context of cross-validation, is to use a criterion which is the average of the cross-validatory criteria corresponding to the different curves. For each of these methods, “average” might be defined in a weighted sense, where the weights represent the respective subsample sizes.

**2.5. Bootstrap calibration**

Bootstrap calibration is along conventional lines, as follows. Having constructed smoothed estimators  $\widehat{X}_i$  and  $\widehat{Y}_i$  of the functions  $X_i$  and  $Y_i$ , respectively, pool them into the class

$$\mathcal{Z} = \{Z_1, \dots, Z_{m+n}\} = \{\widehat{X}_1, \dots, \widehat{X}_m\} \cup \{\widehat{Y}_1, \dots, \widehat{Y}_n\}. \tag{2.5}$$

By sampling randomly, with replacement, from  $\mathcal{Z}$ , derive two independent re-samples  $\{X_1^*, \dots, X_m^*\}$  and  $\{Y_1^*, \dots, Y_n^*\}$ ; compute

$$\bar{F}_X^*(z) = \frac{1}{m} \sum_{i=1}^m I(X_i^* \leq z), \quad \bar{F}_Y^*(z) = \frac{1}{n} \sum_{i=1}^n I(Y_i^* \leq z);$$

and finally, calculate

$$\bar{T}^* = \int \{\bar{F}_X^*(z) - \bar{F}_Y^*(z)\}^2 \mu(dz). \tag{2.6}$$

Of course,  $\bar{F}_X^*$  and  $\bar{F}_Y^*$  are bootstrap versions of the actual empirical distribution functionals,

$$\bar{F}_X(z) = \frac{1}{m} \sum_{i=1}^m I(X_i \leq z), \quad \bar{F}_Y(z) = \frac{1}{n} \sum_{i=1}^n I(Y_i \leq z), \tag{2.7}$$

which we would have computed if we had access to the full-function data  $X_i$  and  $Y_i$ . Likewise,  $\bar{T}$  is the ideal, but impractical, test statistic that we would have used if we had those data:

$$\bar{T} = \int \{\bar{F}_X(z) - \bar{F}_Y(z)\}^2 \mu(dz). \tag{2.8}$$

Suppose the desired critical level for the test is  $1 - \alpha$ . By repeated resampling from  $\mathcal{Z}$  we can compute, by Monte Carlo means, the critical point,  $\hat{t}_\alpha$  say, given by

$$P(\bar{T}^* \geq \hat{t}_\alpha \mid \mathcal{Z}) = \alpha. \tag{2.9}$$

We reject the null hypothesis  $H_0$ , that there is no difference between the distributions of  $X$  and  $Y$ , if  $\widehat{T} > \widehat{t}_\alpha$ .

At this point, some of the potential difficulties of two-sample hypothesis testing become clear. Regardless of how we smooth the data, the conditional expected value of  $\widehat{F}_X^* - \widehat{F}_Y^*$ , given  $\mathcal{Z}$ , is exactly zero. However, even if  $H_0$  is correct,  $E(\widehat{F}_X - \widehat{F}_Y)$  will generally not vanish, owing to the different natures of the datasets providing information about the functions  $X_i$  and  $Y_i$  (for example, different distributions of the sampling times), and the different ways we constructed  $\widehat{X}_i$  and  $\widehat{Y}_i$  from those data. Therefore, the test statistic  $\widehat{T}$  suffers biases which are not reflected in its bootstrap form,  $\overline{T}^*$ , and which can lead to loss of power for the test. Of course, this problem would vanish if we could use  $\widehat{F}_X - \widehat{F}_Y$  in place of  $\widehat{F}_X^* - \widehat{F}_Y^*$  for the test; that is, if we could employ  $\overline{T}$  instead of  $\widehat{T}$ . But in practice, that is seldom possible.

**2.6. Choice of the measure  $\mu$**

Recall from (2.4) that our application of the measure  $\mu$ , to calculate the statistic  $\widehat{T}$ , proceeds by simulating the stochastic process  $M$  of which  $\mu$  defines the distribution. Therefore it is necessary only to construct  $M$ . It is satisfactory to define  $M$  in terms of its Karhunen-Loève expansion,

$$M(t) = \gamma(t) + \sum_{i=1}^{\infty} \zeta_i \phi_i(t),$$

where  $\gamma(t) = E\{M(t)\}$ , the functions  $\phi_i$  form a complete orthonormal sequence on  $\mathcal{I}$ , and the random variables  $\zeta_i$  are uncorrelated and have zero mean.

We take  $\zeta_i = \theta_i \xi_i$ , where  $\theta_1 > \theta_2 > \dots > 0$  are positive constants decreasing to zero, and the random variables  $\xi_i$  are independent and identically distributed with zero means and unit variances. The functions  $\phi_i$  could be chosen to be the orthonormal functions corresponding to a principal component analysis of the dataset  $\mathcal{Z}$  at (2.5). In this case they would form the sequence of eigenvectors of a linear operator, the kernel of which is the function

$$L(s, t) = \frac{1}{m+n} \sum_{i=1}^{m+n} \{Z_i(s) - \overline{Z}(s)\} \{Z_i(t) - \overline{Z}(t)\},$$

where  $\overline{Z} = (m+n)^{-1} \sum_{i \leq m+n} Z_i$ . The constants  $\theta_i$  would in that case be given by the square-root of the corresponding eigenvalues.

However, exposition is simpler if we take the  $\phi_i$ 's to be a familiar orthonormal sequence, such as the cosine sequence on  $\mathcal{I}$ :

$$\phi_1 \equiv 1, \quad \phi_{i+1}(x) = 2^{\frac{1}{2}} \cos(i\pi x) \quad \text{for } i \geq 1.$$



(Recall that  $\mathcal{I} = [0, 1]$ .) In particular, this makes it easier to describe the smoothness of the functions  $M$ . If  $\theta_i = O(i^{-a})$  as  $i \rightarrow \infty$ , where  $a > 3$ , and if the common distribution of the  $\xi_i$ 's is compactly supported, then there exist  $C, c > 0$  such that, with probability 1,  $|M''(s) - M''(t)| \leq C |s - t|^c$ . This is the level of regularity that our theoretical properties require of the distribution of  $M$ . Numerical work in Section 4 argues that in practice it is adequate to take the process  $M$  to have a light-tailed distribution, such as the Gaussian; it is not necessary to assume the distribution is compactly supported.

While choice of the basis functions  $\phi_i$  has little influence on level, it does affect power. In the setting of power against local alternatives, which we discuss in Section 3.4, a theoretically optimal choice of the basis depends on the nature of the local alternative. However, the latter cannot be estimated consistently from data, and so in empirical work (see Section 4), we take an approach which is familiar in more conventional problems: we choose the basis using principal-component methods.

### 3. Theoretical Properties

#### 3.1. Overview

Section 3.2 gives a simple approximation,  $\tilde{T}$ , to  $\hat{T}$ ; Section 3.3 treats a centering term,  $D$ , which represents the main difference between  $\tilde{T}$  and  $\bar{T}$ ; and Section 3.4 describes asymmetric properties of  $\bar{T}$ . These steps in our argument culminate in Section 3.5, which draws our main conclusions. In particular, Section 3.5 combines results of Sections 3.2–3.4 to give conditions under which the statistics  $\hat{T}$  and  $\bar{T}$  have the same asymptotic properties. It also shows that the practical statistic  $\hat{T}$  leads to tests with the same asymptotic level as its idealised counterpart  $\bar{T}$ , and to the same asymptotic power against local alternatives.

#### 3.2. Main approximation property

Let  $\bar{F}_X$ ,  $\bar{F}_Y$  and  $\bar{T}$  be the quantities defined at (2.7) and (2.8). In Section 2.5 we discussed the fact that the bootstrap form of  $\bar{T}$  does not adequately reflect differences between the functionals  $\hat{F}_X$  and  $\hat{F}_Y$  on which the practicable test statistic  $\hat{T}$  is based. Our main result in the present section shows that the main aspects of these potential problems can be encapsulated quite simply in terms of a difference between expected values.

In particular, under very mild conditions, difficulties associated with stochastic variability of  $\hat{F}_X - \hat{F}_Y$  are negligibly small; and the impact of the difference,

$$D(z) = E\{\bar{F}_X(z) - \bar{F}_Y(z)\} - E\{\hat{F}_X(z) - \hat{F}_Y(z)\},$$

can be summarised very simply. Theorem 1 below shows that  $\widehat{T}$  is closely approximated by

$$\widetilde{T} = \int \{ \bar{F}_X(z) - \bar{F}_Y(z) - D(z) \}^2 \mu(dz). \quad (3.1)$$

The sets of assumptions A.1 and A.2, used for Theorems 1 and 2 respectively, will be collected together in Section 5. Proofs of the results are given in Hall and Van Keilegom (2006).

**Theorem 1.** *If the measure  $\mu$  has no atoms, and assumption A.1 holds, then*

$$|\widehat{T}^{\frac{1}{2}} - \widetilde{T}^{\frac{1}{2}}| = o_p(m^{-\frac{1}{2}} + n^{-\frac{1}{2}}). \quad (3.2)$$

To interpret this result, note that, under the null hypothesis the distributions of  $X$  and  $Y$  are identical,  $\bar{T}$  is of size  $m^{-1} + n^{-1}$ . (Further discussion of this point is given in Section 3.4.) The quantity,  $D(z)$ , can only increase this size. Result (3.2) asserts that the difference between  $\widehat{T}^{1/2}$  and  $\widetilde{T}^{1/2}$  is actually of smaller order than the asymptotic sizes of either  $\widehat{T}^{1/2}$  or  $\widetilde{T}^{1/2}$ , and so  $D(z)$  captures all the main issues that will affect the power and level accuracy of the statistic  $\widehat{T}$ , compared with those of  $\bar{T}$ .

### 3.3. Properties of $D(z)$

First we summarise properties of  $D(z)$  when  $(2r - 1)$ st degree local-polynomial estimators are employed. Using arguments similar to those in Hall and Van Keilegom (2006), it may be shown that for functions  $z$  that are sufficiently smooth, and for each  $\eta > 0$ ,

$$P(\widehat{X}_i \leq z) = P(X_i \leq z) + O\{h_{X_i}^{2r-\eta} + (m_i h_{X_i})^{\eta-1}\}. \quad (3.3)$$

This result, and its analogue for the processes  $\widehat{Y}_i$  and  $Y_i$ , lead to the following result. Under the null hypothesis, and for each  $\eta > 0$ ,

$$D(z) = O\left[ \frac{1}{m^{1-\eta}} \sum_{i=1}^m \{h_{X_i}^{2r} + (m_i h_{X_i})^{-1}\} + \frac{1}{n^{1-\eta}} \sum_{j=1}^n \{h_{Y_j}^{2r} + (n_j h_{Y_j})^{-1}\} \right]. \quad (3.4)$$

Neglecting the effects of  $\eta$ , assuming that the subsample sizes  $m_i$  and  $n_i$  are close to a common value  $\nu$ , say, and supposing that the bandwidths  $h_{X_i}$  and  $h_{Y_j}$  are also taken of similar sizes, (3.4) suggests allowing those bandwidths to be of size  $\nu^{-1/(2r+1)}$ . Then  $D(z) = O(\nu^{-2r/(2r+1)})$ .

While this approach is of interest, the extent of reduction in subsampling effects under  $H_0$  can often be bettered by taking the bandwidths  $h = h_{X_i} = h_{Y_j}$  to be identical, all  $1 \leq i \leq m$  and  $1 \leq j \leq n$ . This allows the quantities that contribute the dominant bias terms, involving  $h_{X_i}^{2r}$  and  $h_{Y_j}^{2r}$  in (3.3) and

their analogues for the  $Y$ -sample, to cancel perfectly. That reduces the bias contribution, from the  $2r$ th to the  $2(r + 1)$ st power of bandwidth.

Using identical bandwidths makes local-linear methods, which correspond to taking  $r = 1$  in the formulae above, particularly attractive for at least two reasons. First, the contribution of bias is reduced to that which would arise through fitting third-degree, rather than first-degree, polynomials in the case of non-identical bandwidths, yet the greater robustness of first-degree fitting is retained. Secondly, the appropriate bandwidth is now close to  $\nu^{-1/5}$ , the conventional bandwidth size for estimating the functions  $X_i$  and  $Y_i$  as functions in their own right. This suggests that tried-and-tested bandwidth selectors, such as those discussed in Section 2.4, could be used.

The mathematical property behind the common-bandwidth recommendation is the following more detailed version of (3.3), which for simplicity we give only in the local-linear case, i.e.,  $r = 1$ . If  $h = h_{X_i}$  for each  $i$  and  $h$  is of size  $\nu^{-1/5}$ , or larger, then for each  $\eta > 0$ ,

$$P(\widehat{X}_i \leq z) = P(X_i + \frac{1}{2} \kappa_2 h^2 X_i'' \leq z) + O\{h^{4-\eta} + (\nu h)^{\eta-1}\} \tag{3.5}$$

uniformly in smooth functions  $z$ , where  $\kappa_2 = \int u^2 K(u) du$ . If  $H_0$  holds, and we use the bandwidth  $h$  for the  $Y$  data as well as for the  $X$  data, then

$$P(X_i + \frac{1}{2} \kappa_2 h^2 X_i'' \leq z) = P(Y_i + \frac{1}{2} \kappa_2 h^2 Y_i'' \leq z) \tag{3.6}$$

for each  $i$  and each function  $z$ . Therefore (3.5) implies that, under  $H_0$  and assuming that the subsample sizes are all similar to  $\nu$ ,

$$P(\widehat{X}_i \leq z) - P(\widehat{Y}_j \leq z) = O\{h^{4-\eta} + (\nu h)^{\eta-1}\}$$

for  $1 \leq i \leq m$  and  $1 \leq j \leq n$ .

Optimising the right-hand side of (3.6) with respect to  $h$  suggests using a relatively conventional bandwidth selector of size  $\nu^{-1/5}$ . A small value of the right-hand side of (3.6) implies that  $D$  is close to zero, which in turn ensures that  $\widehat{T}$  (which is close to  $\widetilde{T}$ , as shown in Section 3.2) is well-approximated by  $\overline{T}$ . Such a result implies that the bootstrap test is unlikely to reject  $H_0$  simply because of poor choice of smoothing parameters; see Section 5.2 for further discussion.

We conclude with a concise statement of (3.5). Given sequences  $a_m$  and  $b_m$  of positive numbers, write  $a_m \asymp b_m$  to denote that  $a_m/b_m$  is bounded away from zero and infinity as  $n \rightarrow \infty$ . The reader is referred to Section 5.2 for a statement of assumption A.2 for Theorem 2. It includes the condition that all bandwidths  $h = h_{X_i}$  are identical, and  $h \asymp \nu^{-1/q}$  where  $3 < q < \infty$ .

**Theorem 2.** *If  $r = 1$  and A.2 holds then for each  $\eta > 0$ ,*

$$P(\widehat{X}_i \leq z) - P(X_i + \frac{1}{2} \kappa_2 h^2 X_i'' \leq z) = \begin{cases} O(\nu^{\eta - \frac{(q-1)^2}{4q}}) & \text{if } 3 < q \leq 5 \\ O(\nu^{\eta - \frac{4}{q}}) & \text{if } q > 5, \end{cases} \quad (3.7)$$

where the “big oh” terms are of the stated orders uniformly in functions  $z$  that have two Hölder-continuous derivatives, and for  $1 \leq i \leq m$ .

### 3.4. Asymptotic distribution of $\overline{T}$ , and power

If  $m$  and  $n$  vary in such a way that

$$\frac{m}{n} \rightarrow \rho \in (0, \infty) \quad \text{as } m \rightarrow \infty, \quad (3.8)$$

and if, as prescribed by  $H_0$ , the distributions of  $X$  and  $Y$  are identical, then  $\overline{T}$  satisfies

$$m\overline{T} \rightarrow \zeta \equiv \int \{\zeta_X(z) - \rho^{\frac{1}{2}} \zeta_Y(z)\}^2 \mu(dz). \quad (3.9)$$

Here the convergence is in distribution, and  $\zeta_X(z)$  and  $\zeta_Y(z)$  are independent Gaussian processes with zero means and the same covariance structures as the indicator processes  $I(X \leq z)$  and  $I(Y \leq z)$ , respectively. In particular, the covariance of  $\zeta_X(z_1)$  and  $\zeta_X(z_2)$  is  $F_X(z_1 \wedge z_2) - F_X(z_1)F_X(z_2)$ .

It follows directly from (3.9) that the asymptotic value of the critical point for an  $\alpha$ -level test of  $H_0$ , based on  $\overline{T}$ , is the quantity  $u_\alpha$  such that  $P(\zeta > u_\alpha) = \alpha$ . Analogously, the critical point  $\hat{t}_\alpha$  for the bootstrap statistic  $\overline{T}^*$  (see (2.6) and (2.9)) converges, after an obvious rescaling, to  $u_\alpha$  as sample size increases: under  $H_0$ ,

$$P(m\overline{T} > u_\alpha) \rightarrow \alpha, \quad m\hat{t}_\alpha \rightarrow u_\alpha, \quad (3.10)$$

where the second convergence is in probability. Of course, these are conventional properties of bootstrap approximations. In Section 3.5 we discuss conditions that are sufficient for the practical test statistic  $\widehat{T}$ , rather than its ideal form  $\overline{T}$ , to have asymptotically correct level; see (3.17).

Power properties under local alternatives are also readily derived. In particular, if  $F_Y$  is fixed and

$$F_X(z) = F_Y(z) + m^{-\frac{1}{2}} c \delta(z), \quad (3.11)$$

where  $\delta$  is a fixed function and  $c$  is a constant, then with convergence interpreted in distribution,

$$m\overline{T} \rightarrow \int \{\zeta_{Y1}(z) + c\delta(z) - \rho^{\frac{1}{2}} \zeta_{Y2}(z)\}^2 \mu(dz), \quad (3.12)$$

of which (3.9) is a special case. In (3.12),  $\zeta_{Y1}$  and  $\zeta_{Y2}$  are independent Gaussian processes each with zero mean and the covariance structure of  $\zeta_Y$ .

From this result and the second part of (3.10) it is immediate that, provided  $\delta$  is not almost surely zero with respect to  $\mu$  measure, a test based on the ideal statistic  $\bar{T}$ , but using the bootstrap critical point  $\hat{t}_\alpha$ , is able to detect departures proportional to  $m^{-1/2} \delta$ :

$$\lim_{c \rightarrow \infty} \liminf_{m \rightarrow \infty} P_c(\bar{T} > \hat{t}_\alpha) = 1, \tag{3.13}$$

where  $P_c$  denotes probability under the model where  $F_Y$  is fixed and  $F_X$  is given by (3.11). In Section 3.5 we note that if we use a common bandwidth, and if the subsample sizes are not too much smaller than the sample sizes  $m$  and  $n$ , then the same result holds true for  $\hat{T}$ .

Proofs of (3.9), (3.10) and (3.12) are straightforward. They do not require convergence of function-indexed empirical processes to Gaussian processes, and proceed instead via low-dimensional approximations to those empirical processes.

### 3.5. Sufficient conditions for $\hat{T}$ and $\bar{T}$ to have identical asymptotic distributions under $H_0$

Assume (3.8) and the conditions of Theorem 2 for both the  $X$  and  $Y$  populations, and in particular that all the subsample sizes  $m_i$  and  $n_i$  are of the same order in the sense that

$$\nu \asymp \min_{1 \leq i \leq m} m_i \asymp \max_{1 \leq i \leq m} m_i \asymp \min_{1 \leq i \leq n} n_i \asymp \max_{1 \leq i \leq n} n_i \tag{3.14}$$

as  $m, n \rightarrow \infty$ . Take  $h \asymp \nu^{-1/5}$ . Then Theorem 2, and its analogue for the  $Y$  sample, imply that under  $H_0$ ,

$$D(z) = O(\nu^{\eta - \frac{4}{5}}) \tag{3.15}$$

uniformly in functions  $z$  with two Hölder-continuous derivatives, for each  $\eta > 0$ .

Theorem 1 implies that, in order for the practical statistic  $\hat{T}$ , and its “ideal” version  $\bar{T}$ , to have identical asymptotic distributions, it is necessary only that  $D(z)$  be of smaller order than the stochastic error of  $\bar{F}_X - \bar{F}_Y$ . Equivalently, if (3.8) holds,  $D(z)$  should be of smaller order than  $m^{-1/2}$ , uniformly in functions  $z$  with two Hölder-continuous derivatives. For that to be true it is sufficient, in view of (3.15), that

$$m = O(\nu^{\frac{8}{5} - \eta}) \tag{3.16}$$

for some  $\eta > 0$ .

It follows from (3.10) that, provided (3.16) holds and the null hypothesis is correct,

$$P(m\hat{T} > u_\alpha) \rightarrow \alpha. \tag{3.17}$$

This is the analogue of the first part of (3.10) for  $\widehat{T}$  rather than its ideal form  $\overline{T}$ . Similarly, (3.13) holds for  $\widehat{T}$  rather than  $\overline{T}$ , if a common bandwidth is used and the subsample sizes satisfy (3.14). This confirms the ability of the practicable test, based on  $\widehat{T}$ , to detect semiparametric departures from the null hypothesis.

Condition (3.16) is surprisingly mild. It asserts that, in order for the effects of estimating  $X_i$  and  $Y_i$  to be negligible, it is sufficient that the subsample sizes  $m_i$  and  $n_i$  be of larger order than the 5/8th root of the smaller of the two sample sizes,  $m$  and  $n$ .

#### 4. Numerical Properties

Suppose that  $S_{ij}$  ( $1 \leq i \leq m; 1 \leq j \leq m_i$ ) are i.i.d. uniform on  $[0, 1]$  and that  $T_{ij}$  ( $1 \leq i \leq n; 1 \leq j \leq n_i$ ) are i.i.d. with density  $2 - b + 2(b - 1)t$  for  $0 \leq t \leq 1$  and  $0 < b < 2$ . Note that this last density reduces to the uniform density when  $b = 1$ . We take  $b = 1.2$ . The errors  $\delta_{ij}$  and  $\epsilon_{ij}$  are independent and come from a normal distribution with mean zero and standard deviation  $\sigma = 0.1$  and  $0.3$ , respectively. Suppose that  $X_1, \dots, X_m$  are identically distributed as  $X$ , where  $X(t) = \sum_{k \geq 1} c_k N_{kX} \psi_k(t)$ ,  $c_k = e^{-k/2}$ ,  $N_{kX}$  ( $k \geq 1$ ) are i.i.d. standard normal random variables and  $\psi_k(t) = 2^{1/2} \sin\{(k - 1)\pi t\}$  ( $k > 1$ ) and  $\psi_1 \equiv 1$  are orthonormal basis functions. Similarly, let  $Y_1, \dots, Y_n$  be identically distributed as  $Y$ , where

$$Y(t) = \sum_{k=1}^{\infty} c_k N_{kY1} \psi_k(t) + a \sum_{k=1}^{\infty} a_k N_{kY2} \psi_k^*(t).$$

Here  $N_{kY1}$  and  $N_{kY2}$  are i.i.d. standard normal variables,  $a \geq 0$  controls the deviation from the null model ( $a = 0$  under  $H_0$ ),  $a_k = k^{-2}$ , and

$$\psi_k^*(t) = \begin{cases} 1 & \text{if } k = 1 \\ 2^{1/2} \sin\{(k - 1)\pi(2t - 1)\} & \text{if } k \text{ is odd and } k > 1 \\ 2^{1/2} \cos\{(k - 1)\pi(2t - 1)\} & \text{if } k \text{ is even} \end{cases}$$

are orthonormal basis functions. For practical reasons, we truncate the infinite sum in the definition of  $X(t)$  and  $Y(t)$  at  $k = 15$ . Define  $U_{ij} = X_i(S_{ij}) + \delta_{ij}$  ( $1 \leq i \leq m; 1 \leq j \leq m_i$ ) and  $V_{ij} = Y_i(T_{ij}) + \epsilon_{ij}$  ( $1 \leq i \leq n; 1 \leq j \leq n_i$ ). Finally,  $M_1, \dots, M_N$  are independent and have the same distribution as  $M$ , where  $M(t) = \sum_{k \geq 1} c_k N_{kZ} \phi_k(t)$ ,  $N_{kZ}$  are i.i.d. standard normal variables, and  $\phi_k(t) = 2^{1/2} \cos\{(k - 1)\pi t\}$  ( $k > 1$ ) and  $\phi_1 \equiv 1$  are orthonormal functions. We take  $N = 50$  and truncate the infinite sum after 15 terms. The simulation results are based on 500 samples, and the critical values of the test are obtained from 250 bootstrap samples. The functions  $X(t)$  and  $Y(t)$  are estimated by means of

local-linear smoothing. The bandwidth is chosen to minimise the cross-validation type criterion

$$\frac{1}{m_i m} \sum_{i=1}^m \sum_{j=1}^{m_i} \{U_{ij} - \widehat{X}_i^{-(j)}(S_{ij})\}^2 + \frac{1}{n_i n} \sum_{i=1}^n \sum_{j=1}^{n_i} \{V_{ij} - \widehat{Y}_i^{-(j)}(T_{ij})\}^2,$$

where  $\widehat{X}_i^{-(j)}$  is the estimated regression curve without using observation  $j$  (and similarly for  $\widehat{Y}_i^{-(j)}$ ). The function  $K$  is the biquadratic kernel,  $K(u) = (15/16)(1-u^2)^2 I(|u| \leq 1)$ .

The results for  $m = n = 15, 25, 50$  and  $m_1 = n_1 = 20$  and  $100$  are summarised in Figure 1. The level of significance is  $\alpha = 0.05$  and is indicated in the figure. The graphs show that under the null hypothesis the level is well respected and that the power increases for larger values of  $m, n, m_1, n_1$  and  $a$ . The value of the subsample sizes  $m_1$  and  $n_1$  has limited impact on the power, whereas this is clearly not the case for the sample sizes  $m$  and  $n$ . Other settings that are not reported here (equal variances in the two populations, bigger sample and subsample sizes, ...) show similar behavior for the power curves.

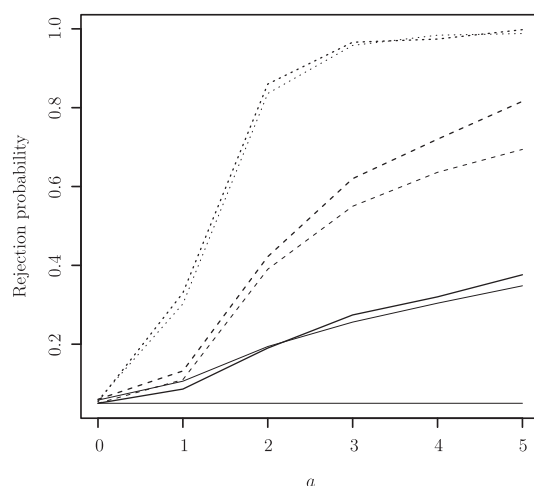


Figure 1. Rejection probabilities for  $m = n = 15$  (full curve),  $m = n = 25$  (dashed curve) and  $m = n = 50$  (dotted curve). The thin curves correspond to  $m_1 = n_1 = 20$ , the thick curves to  $m_1 = n_1 = 100$ . The null hypothesis holds for  $a = 0$ .

In Section 3.3 we explained why it is recommended to take  $h_{X_i} = h_{Y_j} = h$ . We now verify in a small simulation study that identical bandwidths indeed lead to higher power. Consider the same model as above, except that now the standard

deviations of the errors  $\delta_{ij}$  and  $\epsilon_{ij}$  are 0.2 and 0.5, respectively. Take  $m = n = 15$ , 25, and 50,  $m_1 = 20$  and  $n_1 = 100$ . Figure 2 shows the power curves for this model. The rejection probabilities are obtained using either identical bandwidths (estimated by means of the above cross-validation procedure), or using different bandwidths for each sample (estimated by means of a cross-validation procedure for each sample). The graph suggests that under  $H_0$  the empirical level is close to the nominal level in both cases. The power is, however, considerably lower when different bandwidths are used than when the same bandwidth is used for both samples. The differences in power, for different bandwidths, are not so marked when the values of  $m_1$  and  $n_1$  are similar, but in such cases there is little motivation for taking the bandwidths to be quite different.

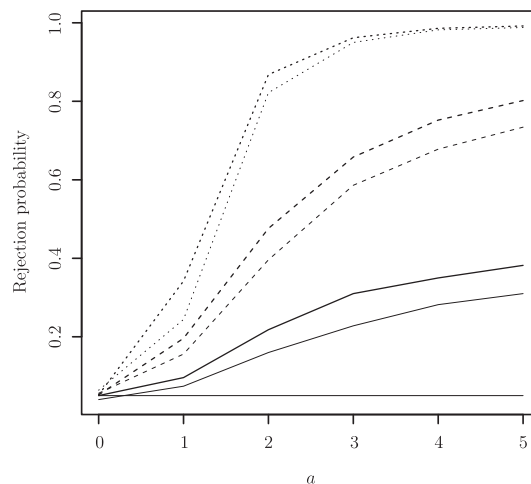


Figure 2. ejection probabilities for  $m = n = 15$  (full curve),  $m = n = 25$  (dashed curve) and  $m = n = 50$  (dotted curve). The thin curves are obtained by using different bandwidths for each sample, the thick curves use the same bandwidth. In all cases,  $m_1 = 20$  and  $n_1 = 100$ . The null hypothesis holds for  $a = 0$ .

Finally, we apply the proposed method to temperature data collected by 224 weather stations across Australia (see Hall and Tajvidi (2002)). The data set consists of monthly average temperatures in degrees Celcius. The starting year of the registration of these data depends on the weather station and varies from 1856 to 1916. Data were registered until 1993. Figure 3 shows the temperature curves of the 224 weather stations during 1990. We restrict attention to the curves between 1914 and 1993, in order to have complete data information for (almost) all weather stations. We are interested in whether weather patterns



in Australia have changed during this time period. To test this hypothesis, we split the time interval in 4 periods of equal length, namely from 1914 to 1933 (period 1), from 1934 to 1953 (period 2), from 1954 to 1973 (period 3) and from 1974 to 1993 (period 4), and execute the proposed pairwise tests on each of the six possible pairs. Curves containing missing monthly temperature values are eliminated from the analysis.

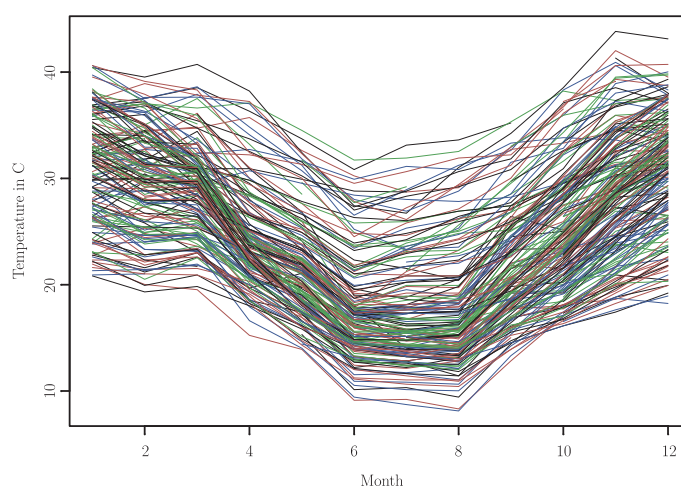


Figure 3. Temperature curves for the 224 weather stations in 1990.

We obtain the optimal bandwidth for each combination of two periods by means of the above cross-validation procedure. In each case we find  $h = 2.2$ . Next, we need to determine an appropriate measure  $\mu$  or, equivalently, an appropriate process  $M$ . For this, we follow the procedure described in Section 2.6: let  $\phi_i(t)$  and  $\theta_i^2$  be the eigenfunctions and eigenvalues corresponding to a principal component analysis (PCA) of the dataset, and truncate the infinite sum at four terms. The variables  $\xi_i$  ( $1 \leq i \leq 4$ ) are taken as independent standard normal variables, and the mean function  $\gamma(t)$  is estimated empirically. The estimation of these functions is carried out by using the PCA routines available on J. Ramsay's homepage (<http://ego.psych.mcgill.ca/misc/fda/>). Next, based on the so-obtained process  $M$ , we calculate the test statistics for each comparison and approximate the corresponding  $p$ -values from 1,000 resamples. The  $p$ -values are summarized in Table 1. At the 5% level, two comparisons (period 1 versus 3 and period 3 versus 4) lead to non-significant  $p$ -values (although borderline), all other  $p$ -values are significant: the data suggest that the average climate in Australia has changed over time.

Table 1.  $P$ -values of the pairwise tests for the Australian weather data.

| Period | Period | $P$ -value |
|--------|--------|------------|
| 1      | 2      | 0.001      |
| 1      | 3      | 0.081      |
| 1      | 4      | 0.024      |
| 2      | 3      | 0.001      |
| 2      | 4      | 0.000      |
| 3      | 4      | 0.053      |

## 5. Assumptions for Section 3

### 5.1. Conditions for Theorem 1

The assumptions are the following, comprising A.1: for some  $\eta > 0$ ,

$$\min \left( \min_{1 \leq i \leq m} m_i, \min_{1 \leq i \leq n} n_i \right) \rightarrow \infty, \quad (5.1)$$

$$\max_{1 \leq i \leq m} h_{X_i} + \max_{1 \leq j \leq n} h_{Y_j} \rightarrow 0, \quad \min_{1 \leq i \leq m} (m_i^{1-\eta} h_{X_i}) + \min_{1 \leq j \leq n} (n_j^{1-\eta} h_{Y_j}) \rightarrow \infty, \quad (5.2)$$

$$K \text{ is a bounded, symmetric, compactly-supported probability density,} \quad (5.3)$$

the observation times  $S_{ij}$  and  $T_{ij}$  are independent random variables, identically distributed for each  $i$ , and with densities that are bounded away from zero uniformly in  $i$  and in population type. (5.4)

Assumption (5.1) asks that the subsample sizes  $m_i$  and  $n_i$  diverge to infinity in a uniform manner as  $m$  and  $n$  grow. This is a particularly mild condition; we do not expect a subsample to provide asymptotically reliable information about the corresponding random functions,  $X_i$  or  $Y_i$ , unless it is large. The first part of (5.2) asks that the bandwidths  $h_{X_i}$  and  $h_{Y_j}$  be uniformly small. Again this is a minimal condition, since bandwidths that converge to zero are necessary for consistent estimation of  $X_i$  and  $Y_i$ . Likewise, the second part of (5.2) is only a little stronger than the assumption that the variances of the estimators of  $X_i$  and  $Y_i$  decrease uniformly to zero.

Assumptions (5.3) and (5.4) are conventional. The latter is tailored to the case of random design, as too is (5.9) below; in the event of regularly spaced design, both can be replaced by simpler conditions.

### 5.2. Conditions for Theorem 2

We need the following notation. Given  $C > 0$  and  $r \in (1, 2]$ , write  $\mathcal{C}_r(C)$  for the class of differentiable functions,  $z$ , on  $\mathcal{I}$  for which: (a)  $\|z'\|_\infty \leq C$ ; (b) if  $1 < r < 2$ ,  $|z'(s) - z'(t)| \leq C |s - t|^{r-1}$  for all  $s, t \in \mathcal{I}$ ; and (c) if  $r = 2$ ,  $z$  has two bounded derivatives and  $\|z''\|_\infty \leq C$ . Given  $d > 0$ , put  $W_d = X + dX''$ , where  $X$  denotes a generic  $X_i$ , and let  $f_{W_d(s)}$  denote the probability density of  $W_d(s)$ .

The assumptions leading to Theorem 2 are the following, comprising A.2:

the kernel  $K$  is a symmetric, compactly-supported probability density with two Hölder-continuous derivatives; (5.5)

$\nu \asymp \min_{1 \leq i \leq m} m_i \asymp \max_{1 \leq i \leq m} m_i$  as  $m \rightarrow \infty$ ; (5.6)

for some  $0 < \eta < 1$  and all sufficiently large  $m$ ,  $m^\eta \leq \nu \leq m^{1/\eta}$ ; (5.7)

the common bandwidth,  $h = h_{X_i}$ , satisfies,

for some  $\eta > 0$ ,  $h \asymp \nu^{-1/q}$  where  $3 < q < \infty$ ; (5.8)

the respective densities  $f_i$  of the observation times  $S_{ij}$  satisfy,

$\sup_{1 \leq i < \infty} \sup_{t \in \mathcal{I}} |f_i''(t)| < \infty$ ,  $\inf_{1 \leq i < \infty} \inf_{t \in \mathcal{I}} f_i(t) > 0$ ; (5.9)

the random function  $X$  has four bounded derivatives, with

$E(|\delta|^s) < \infty$ ,  $E\left\{\sup_{t \in \mathcal{I}} \max_{r=1, \dots, 4} |X^{(r)}(t)|^s\right\} < \infty$  for each  $s > 0$ ; (5.10)

$\sup_{|d| \leq c} \sup_{s \in \mathcal{I}} \|f_{W_d(s)}\|_\infty < \infty$  for some  $c > 0$ ; (5.11)

for  $1 \leq p \leq 2$  and  $c > 0$ , and for each  $C > 0$ ,

$P(W_d \leq z + y) = P(W_d \leq z) + \int_{\mathcal{I}} y(s) P\{W_d \leq z \mid W_d(s) = z(s)\} \times f_{W_d(s)}\{z(s)\} ds + O(\|y\|_\infty^p)$ , (5.12)

uniformly in  $z \in \mathcal{C}_2(C)$ , in  $y \in \mathcal{C}_p(C)$ , and in  $d$  satisfying  $|d| \leq c$ .

Conditions (5.5)–(5.11) are conventional. A proof of (5.12), in the case  $p = 2$  and under the assumption that  $W''$  is well-defined and continuous, is given in Hall and Van Keilegom (2006).

**Acknowledgements**

We are grateful to two reviewers for helpful comments. The research of Van Keilegom was supported by IAP research network grant nr. P5/24 of the Belgian government (Belgian Science Policy).

**References**

Anderson, N. H., Hall, P. and Titterington, D. M. (1994). Two-sample test statistics for measuring discrepancies between two multivariate probability density functions using kernel-based density estimates. *J. Multivariate Anal.* **50**, 41-54.  
 Burdett, K. and Mortensen, D. T. (1998). Wage differentials, employer size, and unemployment. *Internat. Econ. Rev.* **39**, 257-273.

- Cao, R. and Van Keilegom, I. (2006). Empirical likelihood tests for two-sample problems via nonparametric density estimation. *Canad. J. Statist.*, **34**, 61-77.
- Claeskens, G., Jing, B.-Y., Peng, L. and Zhou, W. (2003). Empirical likelihood confidence regions for comparison distributions and ROC curves. *Canad. J. Statist.* **31**, 173-190.
- Cuevas, A., Febrero, M. and Fraiman, R. (2004). An anova test for functional data. *Comput. Statist. Data Anal.* **47**, 111-122.
- Dette, H. and Neumeier, N. (2001). Nonparametric analysis of covariance. *Ann. Statist.* **29**, 1361-1400.
- Fan, J. (1993). Local linear regression smoothers and their minimax efficiencies. *Ann. Statist.* **21**, 196-216.
- Fan, J. and Gijbels, I. (1996). *Local Polynomial Modelling and its Applications*. Chapman and Hall, London.
- Fan, J. and Lin, S.-K. (1998). Test of significance when data are curves. *J. Amer. Statist. Assoc.* **93**, 1007-1021.
- Fan, Y. Q. (1994). Testing the goodness-of-fit of a parametric density-function by kernel method. *Econom. Theory* **10**, 316-356.
- Fan, Y. Q. (1998). Goodness-of-fit tests based on kernel density estimators with fixed smoothing parameter. *Econom. Theory* **14**, 604-621.
- Fan, Y. Q. and Ullah, A. (1999). On goodness-of-fit tests for weakly dependent processes using kernel method. *J. Nonparametr. Statist.* **11**, 337-360.
- Hall, P. and Sheather, S. J. (1988). On the distribution of a Studentized quantile. *J. Roy. Statist. Soc. Ser. B* **50**, 381-391.
- Hall, P. and Tajvidi, N. (2002). Permutation tests for equality of distributions in high-dimensional settings. *Biometrika* **89**, 359-374.
- Hall, P. and Van Keilegom, I. (2006). Two-sample tests in functional data analysis from discrete data. <http://www3.stat.sinica.edu.tw/statistica>.
- Ingster, Y. I. (1993). Asymptotically minimax hypothesis testing for nonparametric alternatives. I, II. *Math. Methods Statist.* **2**, 85-114, 171-189.
- Li, Q. (1999). Nonparametric testing the similarity of two unknown density functions: Local power and bootstrap analysis. *J. Nonparametr. Statist.* **11**, 189-213.
- Koenker, R. and Machado, J. A. F. (1999). Goodness of fit and related inference processes for quantile regression. *J. Amer. Statist. Assoc.* **94**, 1296-1310.
- Koenker, R. and Xiao, Z. J. (2002). Inference on the quantile regression process. *Econometrica* **70**, 1583-1612.
- Locantore, N., Marron, J. S., Simpson, D. G., Tripoli, N., Zhang, J. T. and Cohn, K. L. (1999). Robust principal component analysis for functional data (with discussion). *Test* **8**, 1-73.
- Louani, D. (2000). Exact Bahadur efficiencies for two-sample statistics in functional density estimation. *Statist. Decisions* **18**, 389-412.
- Mortensen, D. T. (1990). Equilibrium wage distributions: A synthesis. In *Panel Data and Labor Market Studies* (Edited by Joop Hartog, Geert Ridder and Jules Theeuwé), 279-96. North-Holland, Amsterdam.
- Ramsay, J. O. and Silverman, B. W. (1997). *Functional Data Analysis*. Springer, New York.
- Ramsay, J. O. and Silverman, B. W. (2002). *Applied Functional Data Analysis: Methods and Case Studies*. Springer, New York.

- Ruppert, D. and Wand, M. P. (1994). Multivariate locally weighted least squares regression. *Ann. Statist.* **22**, 1346-1370.
- Shen, Q. and Faraway, J. (2004). An  $F$  test for linear models with functional responses. *Statist. Sinica* **14**, 1239-1257.
- Spitzner, D. J., Marron, J. S. and Essick, G. K. (2003). Mixed-model functional ANOVA for studying human tactile perception. *J. Amer. Statist. Assoc.* **98**, 263-272.

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

E-mail: p.hall@stat.unimelb.edu.au

Institut de Statistique, Université catholique de Louvain, Voie du Roman Pays 20, B-1348 Louvain-la-Neuve, Belgium.

E-mail: vankeilegom@stat.ucl.ac.be

(Received August 2005; accepted March 2006)