

A spline approach to nonparametric test of hypothesis

Ronaldo Dias and Nancy L. Garcia

Universidade Estadual de Campinas

Abstract: We propose a test of hypothesis for the closeness of two distributions whose test statistic is asymptotically normal. The divergent is based on the estimation procedure developed in Dias (2000) using a proxy of symmetrized Kullback-Leibler distance. Simulation results show that for mixture of normal distributions this test is more powerful than Kolmogorov-Smirnov test. As an application we compare acoustic data from several languages in order to identify rhythmic classes.

Key words: Asymptotic theory; bootstrap; B-splines; Kolmogorov-Smirnov test; Kullback-Leibler divergent.

1 Introduction

There are several situations where we have independent samples and wish to test whether they come from the same distribution. If it is possible to conjecture a parametric family for the distribution, life is much easier and parametric tests can be used. However, most of the time we cannot fit a parametric model and a nonparametric test is necessary. (See for example Fan, 1998, and Li, 1996). The same duality appears in estimation problems. Dias (2000) proposed a nonparametric estimator for densities based on a proxy of the symmetrized Kullback-Leibler distance which is consistent (Section 2). Based on this estimator, in Section 3 we propose a test statistic (henceforth called SKL test) which is asymptotically normal. Simulation results show that for mixture of normal distributions SKL test is more powerful than Kolmogorov-Smirnov (K-S) test (Section 4). Also, the normal approximation is achieved even for small samples when the underlying distribution is normal.

As an application, in Section 5, we present an example that comes from linguistic and deals with clustering the natural languages into rhythmic classes. In the linguistic literature it has been conjectured that natural languages are divided into rhythmic classes (cf. Abercrombie, 1967; Pike, 1945, among others). During half a century no reliable phonetic evidence was presented to support this claim. Recently Ramus, Nespov and Mehler (1999), gave evidence that simple statistical properties of the speech signal could discriminate between different rhythmic classes. They analyzed the acoustic signal of 20 sentences of each of the following languages: English, Polish, Dutch, Catalan, Spanish, Italian, French and Japanese. They computed for each sentence the standard deviation of the consonantal intervals (ΔC) and the proportion of time spent in vocalic intervals ($\%V$)

and found that based on these statistics the languages appear to cluster into three groups which correspond precisely to the intuitive notion of rhythmic classes: English, Polish and Dutch represent the accentual class, French, Spanish, Catalan and Italian represent the syllabic class and Japanese represents the moraic class. In their work there is no study for Portuguese. In Section 5, we apply the proposed nonparametric test to some of these languages and find that there is no significant evidence of difference between European and Brazilian Portuguese, English and Dutch and English and European Portuguese, while there is significant difference between Brazilian Portuguese and Catalan and English and Japanese.

2 Previous results

Suppose we have two independent random samples $\mathbf{X} = (X_1, X_2, \dots, X_{n_1})$ with distribution F and $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{n_2})$ with distribution G and we would like to test whether $F = G$. First assume that both F and G are absolutely continuous cumulative distribution functions with $F \ll \mu$ and $G \ll \mu$ for a Lebesgue dominant measure μ . Moreover, assume that $f = \frac{dF}{d\mu}$ and $g = \frac{dG}{d\mu}$, the respective densities of F and G , have compact support \mathcal{X} . Define \mathcal{F}_μ be the class of density functions such that,

$$\mathcal{F}_\mu = \{h : \mathbb{R} \rightarrow [0, \infty) : h(x) = \frac{e^{S(x)}}{\int_{\mathcal{X}} e^{S(x)} d\mu(x)} \text{ and } \int_{\mathcal{X}} e^{S(x)} d\mu(x) < \infty\},$$

where the function S is of the class $C^2(\mathbb{R})$. It is easy to see that the elements in \mathcal{F}_μ are not identifiable since for any function S_1 such that $S_1 = S + c$, we have $e^{S_1}/(\int e^{S_1}) = e^S/(\int e^S)$. We are going to require, as Dias (1998a), that $\int_{\mathcal{X}} S = 0$, to ensure uniqueness of the elements in \mathcal{F}_μ .

Let h be a density with respect to μ . Consider the problem of finding the maximum likelihood estimator of h . It is well known (see for example, Silverman, 1986; Pagan and Ullah, 1999, and Dias, 1994) that such optimization problem is unbounded over the class of all smooth functions. In fact, the optimizer is a sum of delta functions. To avoid the *Dirac's disaster* one might want to apply penalized likelihood procedure or one may assume that h can be well approximated by a function belonging to a finite dimensional space \mathcal{H}_K which is spanned by K (fixed) basis functions, such as Fourier expansion, wavelets, B-splines, natural splines. See, for example, Silverman (1986), Kooperberg and Stone (1991), Vidakovic (1999), Dias (1998) and Dias (2000). Although this fact might lead one to think that the nonparametric problem becomes a parametric problem, one notices that the number of coefficients can be as large as the number of observations, and there may be difficulties in estimating the density. Moreover, if the number of observations is large, the system of equations for exact solution is too expensive to solve. This is an inheritance from the approximation theory of functions.

In fact, an element of \mathcal{H}_K can be written as

$$h = \frac{e^{S_h}}{\int e^{S_h}}$$

with

$$S_h = \sum_{j=1}^K \theta_j M_j \quad \text{with} \quad \int e^{S_h} < \infty$$

where M_1, \dots, M_K are normalized basis functions that span \mathcal{H}_K such that $\int M_j = 1$. In order to enforce one-to-one correspondence we restrict $\int S_h = 0$ and then $\sum_{j=1}^K \theta_j = 0$, since $\int M_j = 1$. For any $K > 0$, let $\Theta_0 = \{\theta \in \mathbb{R}^K : \sum_{j=1}^K \theta_j = 0\}$.

The vector of coefficients θ are unknown and need to be determined. One of the most common standard statistical procedure in nonparametric estimation, is to determine θ using maximum likelihood method.

Moreover, the vector of coefficients θ does not have a statistical meaning as percentiles, moments, skewness. However, the set of coefficients (as a whole) is extremely important to determining the shape of the density h . Observe that changing the dimension, e.g. by 1, might change completely the estimation of the coefficients for the new set of basis functions. However, for large K , it does not change significantly the estimate of h . In fact, we use this as a stopping criteria in adaptive procedures.

Assuming that the densities f and g belong to \mathcal{F}_μ , we have that there exist K_1, K_2 such that f and g are well approximated by functions in \mathcal{H}_{K_1} and \mathcal{H}_{K_2} respectively. Consequently, there exist vectors $\theta = (\theta_1, \dots, \theta_{K_1})$, $\psi = (\psi_1, \dots, \psi_{K_2})$ such that the log-likelihood of \mathbf{X} and \mathbf{Y} are given by

$$L_{K_1}(\theta|\mathbf{X}) = \frac{1}{n_1} \sum_{i=1}^{n_1} \langle \theta, M(X_i) \rangle_{K_1} - \log \int e^{\langle \theta, M \rangle_{K_1}} \quad (2.1)$$

and

$$L_{K_2}(\psi|\mathbf{Y}) = \frac{1}{n_2} \sum_{i=1}^{n_2} \langle \psi, M(Y_i) \rangle_{K_2} - \log \int e^{\langle \psi, M \rangle_{K_2}}. \quad (2.2)$$

The next results (Lemma 1, Theorem 1, Lemma 2 and Proposition 1) were proved by Dias (2000) in the case that functions M are the normalized B-splines. We are going to enunciate the results for the (2.1) but obviously the results are also valid for (2.2).

Lemma 1 *For a fixed K_1 , $L_{K_1}(\theta|\mathbf{X})$ is concave in θ . Moreover, $L_{K_1}(\theta|\mathbf{X})$ is strictly concave for $\theta \in \Theta_0$. Hence there exists at most one maximizer on Θ_0 .*

It is not difficult to show that $L_{K_1}(\theta|\mathbf{X})$ is continuous and at least twice differentiable in θ for a fixed K . Thus, restrict to Θ_0 one may guarantee a unique density estimate.

The next theorem shows the relationship between the maximizers $\hat{\theta}$ in Θ and θ^* in Θ_0 .

Theorem 1 *If the vector $\hat{\theta}$ maximizes $L_{K_1}(\theta|\mathbf{X})$ then $\theta^* = \hat{\theta} - \frac{1}{K} \sum_{j=1}^{K_1} \hat{\theta}_j$ maximizes $L_{K_1}(\theta|\mathbf{X})$ subject to $\sum_{j=1}^{K_1} \theta_j = 0$. Moreover, θ^* is unique.*

For fixed K_1 , let $\hat{\theta}_{n_1}^{(K_1)}$ be defined as

$$\hat{\theta}_{n_1}^{(K_1)} = \arg \max_{\theta \in \Theta_0} L_{K_1}(\theta|\mathbf{X}). \quad (2.3)$$

Notice that, in fact,

$$L_{K_1}(\theta|\mathbf{X}) = \langle \theta, \bar{M} \rangle_{K_1} - \log \int e^{\langle \theta, M \rangle_{K_1}},$$

then $\hat{\theta}_{n_1}^{(K_1)}$ is the unique solution of the equation

$$h(\theta, \bar{M}(\mathbf{X})) = 0, \quad (2.4)$$

where $\bar{M}(\mathbf{X})$ is a K -dimensional vector with j -th components given by

$$\frac{1}{n_1} \sum_{i=1}^{n_1} M_j(X_i) = \bar{M}_j, \quad j \in \{1, \dots, K_1\}. \quad (2.5)$$

Since $L_{K_1}(\theta|\mathbf{X})$ is at least twice differentiable we have $\hat{\theta}_{n_1}^{(K_1)}$ as the unique solution of the equation,

$$\frac{\partial L_{K_1}(\theta|\mathbf{X})}{\partial \theta} := h(\theta, M^*(\mathbf{X})) = 0, \quad (2.6)$$

where, $M^* = (1/K) \sum_{j=1}^K \bar{M}_j$ and $h : \Theta_0 \times [0, \infty)^K \rightarrow \mathbb{R}^K$ with j -th entry,

$$h_j(\theta, \mathbf{u}) = u_j - \frac{\int \exp(\langle \theta, M(z) \rangle_{K_1}) M_j(z) dz}{\int \exp(\langle \theta, M(z) \rangle_{K_1}) dz}, \quad (2.7)$$

for $j \in \{1, \dots, K_1\}$. Therefore, $\hat{\theta}_{n_1}^{(K_1)}$ is an M-estimator and since $\theta \mapsto h_\theta$ is continuous we have the following result.

Lemma 2 *Let θ_0 be the unique solution of*

$$h(\theta, \int f(x)M(x)d\mu(x)) = 0$$

in Θ_0 , then for fixed K_1 , $\hat{\theta}_{n_1}^{(K_1)} \rightarrow \theta_0$ almost surely as $n_1 \rightarrow \infty$.

Thus, the density estimate is, for fixed K_1

$$\hat{f}_{K_1} = e^{\hat{S} - \log \int e^{\hat{S}}},$$

where $\hat{S} = \langle \hat{\theta}, M \rangle_{K_1}$ with $\hat{\theta} = \hat{\theta}_{n_1}^{(K_1)}$.

One may notice the density estimate \hat{f}_{K_1} strongly depends on the number of basis functions K_1 which regularizes the optimization problem (2.1). In order to provide an appropriate K_1 , one may want to compute the Kullback-Leibler distance between the true f and the random function \hat{f}_{K_1} .

$$d(f, \hat{f}_{K_1}) = \int (\log f - \log \hat{f}_{K_1}) f \quad (2.8)$$

Of course, we cannot compute $d(f, \hat{f}_{K_1})$ from the data, since it requires the knowledge of f . But theoretically we can investigate this distance for the choice of an optimal K_1 in the sense of minimizing $d(f, \hat{f}_{K_1})$. Then, one may define the best K_1 as

$$\hat{K}_1 = \arg \min_{K \in \{1, \dots, K_{max}\}} d(f, \hat{f}_K),$$

for $K_{max} < n_1$. Observe that, in order to obtain \hat{K}_1 , it is sufficient to minimize

$$D_{n_1}(K) = \int f \log \hat{f}_K.$$

Notice that $D_{n_1}(K)$ is a random function of K and also can be approximate by

$$Z_{n_1}(K) = \frac{1}{n_1} \sum_{i=1}^{n_1} \log \hat{f}_K(X_i).$$

Proposition 1 For any fixed K ,

$$D_{n_1}(K) - Z_{n_1}(K) = \sum_{j=1}^K \hat{\theta}_{n_1 j}^{(K)} \left(\int f(x) M_j(x) d\mu(x) - \frac{1}{n_1} \sum_{i=1}^{n_1} M_j(X_i) \right) \longrightarrow 0 \quad (2.9)$$

$n_1 \longrightarrow \infty$ almost surely.

Lemma 3 For K_1, K_2 fixed the density estimates $\hat{f}_{K_1}(\cdot) = f_{K_1}(\cdot | \hat{\theta}_{n_1})$ and $\hat{g}_{K_2}(\cdot) = g_{K_2}(\cdot | \hat{\psi}_{n_2})$ converge pointwise almost surely (a.s.) to $f_{K_1}(\cdot | \theta)$ and $g_{K_2}(\cdot | \psi)$ respectively as n_1, n_2 go to infinity.

Proof. It is enough to show one of the statements above. For fixed x , it is not difficult to check that the map $\theta \mapsto f_m(x|\theta)$ is a continuous map in $\theta \in \Theta_0$, for any $m \in \mathbb{N}$. By Lemma 2 we have $\hat{\theta}_{n_1} \longrightarrow \theta (\in \Theta_0)$ a.s. and so $f_m(x|\hat{\theta}_{n_1}) \longrightarrow f_m(x|\theta)$ almost surely as $n_1 \longrightarrow \infty$. Notice that the null sets of the a.s. convergence do not depend on x then $f(\cdot | \hat{\theta}_{n_1})$ converges pointwise to $f(\cdot)$ a.s. \square

It is important to emphasize that in developing the theory and proving the theorems we followed the approach of non-stochastic K_1 and K_2 . That is, we assume that the dimensions of the approximant spaces are known in advance. However, in practice, K_1 and K_2 are unknown and an adaptive procedure is suggested in

page 5, by using \hat{K}_1 and \hat{K}_2 (which is obtained analogously). Alternative methods can be used, for example a Bayesian point of view can be used via Reversible Jump MCMC, cf. Dias and Gamerman (2002). An stochastic dimension would pose a more difficult problem to a nonparametric test of hypothesis and it is not addressed in this manuscript and it will be left for future research.

3 Hypothesis testing - SKL test statistic

In this section we propose a statistic to test: $H_0 : f = g$ almost surely μ versus the alternative hypothesis $H_1 : f \neq g$ over a set of positive μ -measure. Since this test statistic is based on the symmetrized Kullback-Leibler distance we will call it SKL test.

First we notice that the parameter space $\Theta_0 = \{\theta \in \mathbb{R}^K : \sum_j^K \theta_j = 0\}$ is not an open set and it is a $(K - 1)$ -dimensional manifold in \mathbb{R}^K . Therefore, in order to have any kind of asymptotic normality results we need to reparametrize the problem to $(\theta_1, \dots, \theta_{K-1}) \in \tilde{\Theta}_0$ such that $(\theta_1, \dots, \theta_{K-1}, -\sum_{i=1}^{K-1} \theta_i) \in \Theta_0$. We will continue to call the parameter θ . In this case, the density will be written as

$$f(x | \theta) = \frac{e^{\langle \theta, \tilde{M}(x) \rangle}}{\int_{\mathcal{X}} e^{\langle \theta, \tilde{M}(x) \rangle}} \in \mathcal{F}_\mu \quad (3.1)$$

where

$$\tilde{M}_j(x) = M_j(x) - M_K(x). \quad (3.2)$$

For fixed K_1 and K_2 , we have as a consequence of Cramér's Theorem (see for example, Ferguson, 1996, p.121) the asymptotic normality of the consistent estimator which solves the likelihood equation. For simplicity we will continue to denote

$$\hat{\theta}_{n_1} = (\hat{\theta}_1, \dots, \hat{\theta}_{K_1-1}) \quad (3.3)$$

where $\hat{\theta}_{n_1}^{(K_1)} = (\hat{\theta}_1, \dots, \hat{\theta}_{K_1})$ is given by (2.3). Define similarly $\hat{\psi}_{n_2}$.

Theorem 2 *The estimators $\hat{\theta}_{n_1}$ and $\hat{\psi}_{n_2}$ are asymptotically normal distributed. More specifically, if θ_0 and ψ_0 are the true parameter values, there exists positive definite matrices Σ_1 and Σ_2 such that*

$$\sqrt{n_1}(\hat{\theta}_{n_1} - \theta_0) \rightarrow N_{K_1-1}(0, \Sigma_1) \quad (3.4)$$

$$\sqrt{n_2}(\hat{\psi}_{n_2} - \psi_0) \rightarrow N_{K_2-1}(0, \Sigma_2) \quad (3.5)$$

as $n_1, n_2 \rightarrow \infty$.

Proof. We are going to prove (3.4), (3.5) is completely analogous. We need to show that all $f \in \mathcal{F}_\mu$ satisfy the regularity conditions. First it is obvious that $\tilde{\Theta}_0$ is an open set and the model is identifiable.

Note that $e^{\langle \theta, \tilde{M}(x) \rangle}$ is of the class $\mathcal{C}^\infty(\tilde{\Theta}_0)$. It is easy to verify that $\int_{\mathcal{X}} e^{\langle \theta, \tilde{M}(x) \rangle} < \infty$, since \mathcal{X} is a compact set and $f \in \mathcal{F}_\mu$. In addition,

$$\frac{\partial}{\partial \theta_i} e^{\langle \theta, \tilde{M}(x) \rangle} = \tilde{M}_i(x) e^{\langle \theta, \tilde{M}(x) \rangle} \quad (3.6)$$

$$\frac{\partial^2}{\partial \theta_i \partial \theta_j} e^{\langle \theta, \tilde{M}(x) \rangle} = \tilde{M}_i(x) \tilde{M}_j(x) e^{\langle \theta, \tilde{M}(x) \rangle} \quad (3.7)$$

which exist and are continuous functions for each $(\theta, x) \in \tilde{\Theta}_0 \times \mathcal{X}$. Thus, $f(x|\theta)$ is of the class $\mathcal{C}^\infty(\tilde{\Theta}_0)$ and the partial derivatives may be passed under the integral sign. Let

$$C(\theta) = \int_{\mathcal{X}} e^{\langle \theta, \tilde{M}(x) \rangle} \quad (3.8)$$

$$C_i(\theta) = \int_{\mathcal{X}} \tilde{M}_i(x) e^{\langle \theta, \tilde{M}(x) \rangle} \quad (3.9)$$

$$C_{ij}(\theta) = \int_{\mathcal{X}} \tilde{M}_i(x) \tilde{M}_j(x) e^{\langle \theta, \tilde{M}(x) \rangle}. \quad (3.10)$$

It is easy to verify that

$$\frac{\partial^2 \log f(x|\theta)}{\partial \theta_i \partial \theta_j} = \frac{C_{ij}(\theta)C(\theta) - C_i(\theta)C_j(\theta)}{C(\theta)^2} \quad (3.11)$$

Since $C(\theta)$, $C_i(\theta)$ and $C_{ij}(\theta)$ are continuous functions of θ , then in a closed neighborhood $\mathcal{N}(\theta_0)$ of the true parameter value θ_0 , we have

$$\begin{aligned} C_* &:= \min_{\theta \in \mathcal{N}(\theta_0)} C(\theta) > 0 \\ \bar{C}_i &:= \max_{\theta \in \mathcal{N}(\theta_0)} C_i(\theta) < \infty \\ \bar{C}_{ij} &:= \max_{\theta \in \mathcal{N}(\theta_0)} C_{ij}(\theta) < \infty. \end{aligned}$$

Thus,

$$\left| \frac{\partial^2 \log f(x|\theta)}{\partial \theta_i \partial \theta_j} \right| \leq \frac{|\bar{C}_{ij}|}{C_*} + \frac{|\bar{C}_i \bar{C}_j|}{C_*^2} =: C(i, j) < \infty. \quad (3.12)$$

In a completely analogous way, the third partial derivatives can be bounded.

Let $\mathcal{I}(\theta)$ the Hessian matrix of $\log f(x|\theta)$ with entries

$$\mathcal{I}_{i,k}(\theta) = -\mathbb{E} \left[\frac{\partial^2 \log f(x|\theta)}{\partial \theta_k \partial \theta_i} \right].$$

Then $\mathcal{I}(\theta)$ is positive definite matrix. To see this, observe that

$$\mathcal{I}_{i,k}(\theta) = \text{Cov}(\tilde{M}_i(X), \tilde{M}_j(X)) = \text{Cov}(M_i(X) - M_K(X), M_j(X) - M_K(X)),$$

which is nonnegative definite. But the collection of functions $\{M_1, \dots, M_K\}$ is a basis for the finite dimensional approximant space (e.g., natural cubic spline space) and so the columns of $\mathcal{I}(\theta)$ are linear independent. Consequently, $\mathcal{I}(\theta)$ is a positive definite matrix. \square

To measure the distance between the two distributions F and G we can use the divergent given by:

$$I_S(F, G) := \int (\log f(x) - \log g(x))f(x)d\mu(x) + \int (\log g(y) - \log f(y))g(y)d\mu(y). \quad (3.13)$$

Define the following estimator for $I_S(F, G)$,

$$I_S(\hat{f}, \hat{g}) = \int (\log \hat{f}_{K_1}(x) - \log \hat{g}_{K_2}(x))dF_{n_1}(x) + \int (\log \hat{g}_{K_2}(y) - \log \hat{f}_{K_1}(y))dG_{n_2}(y) \quad (3.14)$$

where F_{n_1} and G_{n_2} are the empirical distribution of F and G respectively. In fact, if we have the random samples $\mathbf{X} = (X_1, X_2, \dots, X_{n_1})$ with distribution F and $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{n_2})$ with distribution G , this estimator is of the form

$$\begin{aligned} I_S(\hat{f}, \hat{g}) &= \frac{1}{n_1} \left(\sum_{i=1}^{n_1} \log \hat{f}_{K_1}(X_i) - \sum_{i=1}^{n_1} \log \hat{g}_{K_2}(X_i) \right) \\ &+ \frac{1}{n_2} \left(\sum_{i=1}^{n_2} \log \hat{g}_{K_2}(Y_i) - \sum_{i=1}^{n_2} \log \hat{f}_{K_1}(Y_i) \right) \\ &= \hat{I}_{1, n_1}(\mathbf{X}) + \hat{I}_{2, n_2}(\mathbf{Y}) \end{aligned} \quad (3.15)$$

Lemma 4 For fixed K_1 and K_2 , $(\hat{I}_{1, n_1}(\mathbf{X}) + \hat{I}_{2, n_2}(\mathbf{Y})) \longrightarrow I_S$ almost surely when $n_1, n_2 \longrightarrow \infty$.

Proof. Call

$$I_1 = \int (\log f(x) - \log g(x))f(x)d\mu(x) \quad (3.16)$$

and

$$I_2 = \int (\log g(y) - \log f(y))g(y)d\mu(y). \quad (3.17)$$

It is enough to show that $\hat{I}_{1, n_1}(\mathbf{X}) \longrightarrow I_1$ almost surely as $n_1, n_2 \longrightarrow \infty$, the

result for $\hat{I}_{2,n_2}(\mathbf{Y})$ is done similarly. Let

$$D_{n_1}^1(K_1) = \int f(x) \log \hat{f}_{K_1}(x) d\mu(x), \quad (3.18)$$

$$D_{n_1,n_2}^2(K_2) = \int f(x) \log \hat{g}_{K_2}(x) d\mu(x), \quad (3.19)$$

$$Z_{n_1}^1(K_1) = n_1^{-1} \sum_{i=1}^{n_1} \log \hat{f}_{K_1}(X_i), \quad (3.20)$$

$$Z_{n_1,n_2}^2(K_2) = n_1^{-1} \sum_{i=1}^{n_1} \log \hat{g}_{K_2}(X_i). \quad (3.21)$$

Note that, $\hat{I}_{1,n_1}(\mathbf{X}) = Z_{n_1}^1(K_1) - Z_{n_1,n_2}^2(K_2)$. Therefore, by adding and subtracting $D_{n_1}^1(K_1)$ and $D_{n_1,n_2}^2(K_2)$, it is sufficient to show

$$(D_{n_1}^1(K_1) - Z_{n_1}^1(K_1)) + (D_{n_1,n_2}^2(K_2) - Z_{n_1,n_2}^2(K_2)) \longrightarrow 0, \quad (3.22)$$

almost surely as $n_1, n_2 \longrightarrow \infty$. For this, observe that,

$$(D_{n_1}^1(K_1) - Z_{n_1}^1(K_1)) = \sum_{j=1}^{K_1-1} \hat{\theta}_{n_1,j} \left(\int f(x) \tilde{M}_{j,1}(x) d\mu(x) - n_1^{-1} \sum_{i=1}^{n_1} \tilde{M}_{j,1}(X_i) \right)$$

and

$$(D_{n_1,n_2}^2(K_2) - Z_{n_1,n_2}^2(K_2)) = \sum_{l=1}^{K_2-1} \hat{\psi}_{n_2,l} \left(\int f(x) \tilde{M}_{l,2}(x) d\mu(x) - n_1^{-1} \sum_{i=1}^{n_1} \tilde{M}_{l,2}(X_i) \right)$$

Following Dias (2000), we have that $\hat{\theta}_{n_1}$ and $\hat{\psi}_{n_2}$ are the maximum likelihood estimators and $\hat{\theta}_{n_1} \longrightarrow \theta_0$ and $\hat{\psi}_{n_2} \longrightarrow \psi_0$ almost surely as $n_1, n_2 \longrightarrow \infty$, where θ_0 and ψ_0 are the true parameter values. Moreover, by the strong law of large numbers,

$$n_1^{-1} \sum_{i=1}^{n_1} \tilde{M}_{j,1}(X_i) \longrightarrow \int f(x) \tilde{M}_{j,1}(x) d\mu(x)$$

and

$$n_1^{-1} \sum_{i=1}^{n_1} \tilde{M}_{l,2}(X_i) \longrightarrow \int f(x) \tilde{M}_{l,2}(x) d\mu(x)$$

almost surely as $n_1 \longrightarrow \infty$, for $j = 1, \dots, K_1 - 1$ and $l = 1, \dots, K_2 - 1$. \square

Theorem 3 For all $f, g \in \mathcal{F}_\mu$ there exists a positive constant σ_I such that

$$\sqrt{n_1}(\hat{I}_{1,n_1}(\mathbf{X}) - I_1) + \sqrt{n_2}(\hat{I}_{2,n_2}(\mathbf{Y}) - I_2) \longrightarrow N(0, \sigma_I)$$

when $(n_1/n_2) \rightarrow 1$ as $n_1, n_2 \longrightarrow \infty$. Note that under H_0 we have $I_1 = I_2 = 0$ and the result turns to be

$$\sqrt{n_1} \hat{I}_{1,n_1}(\mathbf{X}) + \sqrt{n_2} \hat{I}_{2,n_2}(\mathbf{Y}) \longrightarrow N(0, \sigma_I).$$

Proof. Using the notation (3.18)–(3.21) introduced in Lemma 4, it is easy to see that

$$\sqrt{n_1}(\hat{I}_{1,n_1}(\mathbf{X}) - I_1) \quad (3.23)$$

$$\begin{aligned} &= \sqrt{n_1} \left(D_{n_1}^1(K_1) - \int f(x) \log f(x) d\mu(x) \right) \\ &\quad - \sqrt{n_1} \left(D_{n_1, n_2}^2(K_2) - \int f(x) \log g(x) d\mu(x) \right) \\ &\quad + \sqrt{n_1} (Z_{n_1}^1(K_1) - D_{n_1}^1(K_1)) - \sqrt{n_1} (Z_{n_1, n_2}^1(K_2) - D_{n_1, n_2}^1(K_2)) \\ &= \sqrt{n_1} \left\{ \int f(x) \left[\sum_{j=1}^{K_1-1} (\hat{\theta}_j - \theta_j) \tilde{M}_{j,1}(x) \right] d\mu(x) \right\} \\ &\quad - \sqrt{n_1} \left\{ \int f(x) \left[\sum_{l=1}^{K_2-1} (\hat{\psi}_l - \psi_l) \tilde{M}_{l,1}(x) \right] d\mu(x) \right\} \end{aligned} \quad (3.24)$$

$$\begin{aligned} &\quad - \sqrt{n_1} \left\{ \log \int e^{\langle \hat{\theta}, \tilde{M}_1(x) \rangle} d\mu(x) - \log \int e^{\langle \theta, \tilde{M}_1(x) \rangle} d\mu(x) \right\} \\ &\quad + \sqrt{n_1} \left\{ \log \int e^{\langle \hat{\psi}, \tilde{M}_2(x) \rangle} d\mu(x) - \log \int e^{\langle \psi, \tilde{M}_2(x) \rangle} d\mu(x) \right\} \end{aligned} \quad (3.25)$$

$$\begin{aligned} &\quad + \sqrt{n_1} \left\{ \frac{1}{n_1} \left[\sum_{i=1}^{n_1} \sum_{j=1}^{K_1-1} \hat{\theta}_j \tilde{M}_{j,1}(X_i) - \int \sum_{j=1}^{K_1-1} \hat{\theta}_j \tilde{M}_{j,1}(x) f(x) d\mu(x) \right] \right\} \\ &\quad - \sqrt{n_1} \left\{ \frac{1}{n_1} \left[\sum_{i=1}^{n_1} \sum_{l=1}^{K_2-1} \hat{\psi}_l \tilde{M}_{l,2}(X_i) - \int \sum_{l=1}^{K_2-1} \hat{\psi}_l \tilde{M}_{l,2}(x) f(x) d\mu(x) \right] \right\}. \end{aligned}$$

By applying Theorem 2, the Delta Method and the Central Limit Theorem for i.i.d. random variables we get the desired result. \square

Note that the assumption $n_1/n_2 \rightarrow 1$ can be relaxed to $n_1/n_2 \rightarrow c$, as $n_1, n_2 \rightarrow \infty$. This change would just affect (3.24) and (3.25). In fact, applying the Central Limit Theorem, we get

$$\sqrt{n_2}(\hat{\psi}_l - \psi_l) \rightarrow N(0, \sigma^2)$$

where σ^2 is the asymptotic variance. However, what we need is the asymptotic distribution of $\sqrt{n_1}(\hat{\psi}_l - \psi_l)$ and $\sqrt{n_1}(\log \int e^{\langle \hat{\psi}, \tilde{M}_2(x) \rangle} d\mu(x) - \log \int e^{\langle \psi, \tilde{M}_2(x) \rangle} d\mu(x))$. It is immediate to see that

$$\sqrt{n_1}(\hat{\psi}_l - \psi_l) \rightarrow N(0, c\sigma^2)$$

and applying the Delta Method we obtain the asymptotic distribution for (3.25).

Extensions. This procedure can be extended to test closeness of multivariate distribution functions by using tensor product among the B-spline basis. Also, one might consider the dimension of the approximant spaces (K_1 and K_2) to be unknown and estimated from the data using either an adaptive procedure similar to H-splines (Dias, 1998) or a Bayesian approach similar to the one proposed by Dias and Gamerman (2002) for nonparametric regression.

4 Simulation results

Simulations were run in Athlon machine with double 2 GHz processors and used software R (www.r-project.org). All the asymptotic variances given in Theorem 3 were computed using nonparametric bootstrap method.

In the simulations we used as basis functions the well-known B-splines which have bounded support. Although the procedure supposes that the density is continuous and positive on a compact set \mathcal{X} our simulations include test functions which do not have a compact support, e.g., mixture of normal distributions and gamma distributions. Nevertheless, for practical purposes a density with a infinite domain can be approximate by a density with an appropriate compact support. For example, a normal density ϕ_{μ, σ^2} with mean μ and variance σ^2 do not differ significantly from a density on $[\mu - 5\sigma, \mu + 5\sigma]$ proportional to ϕ_{μ, σ^2} . Similarly, for densities in the gamma family. Estimation of the support of a density is a very difficult problem (Hall, Nussbaum and Stern, 1997) which has not been answered appropriately and it will not be addressed in this work.

In order to assess the range of applicability we performed some simulation for small samples using several known distributions. Figure 1 shows that the normal distribution for the test statistic holds even for samples of size 30 when the underlying true distribution is normal. This result was verified using 1000 nonparametric bootstrap resampling of the original data and 1000 independent replications of the sampling distribution. For non-symmetric distributions such as gamma distributions we have a small skewness to the right.

Moreover, we compare SKL test with Kolmogorov-Smirnov (K-S) test which is the most used nonparametric test for comparing continuous distributions. It is well-known that K-S test presents problems in heavy-tailed distributions (see, Mason and Schuenemeyer, 1983, and Mason and Schuenemeyer, 1992). Therefore, we chose to make this comparison in terms of power using 2000 nonparametric bootstrap samples of mixture of normal distributions. Several mixtures of distributions were used. In this work we present a typical example using as a sampling distribution

$$f(x) = .8\phi((x + .5)/.6) + .2\phi((x - \mu)/.6) \quad (4.1)$$

where ϕ is the standard normal density and μ is the mean of the contaminating distribution. Table 1 and Figure 2 show that SKL is consistently more powerful than K-S in this case.

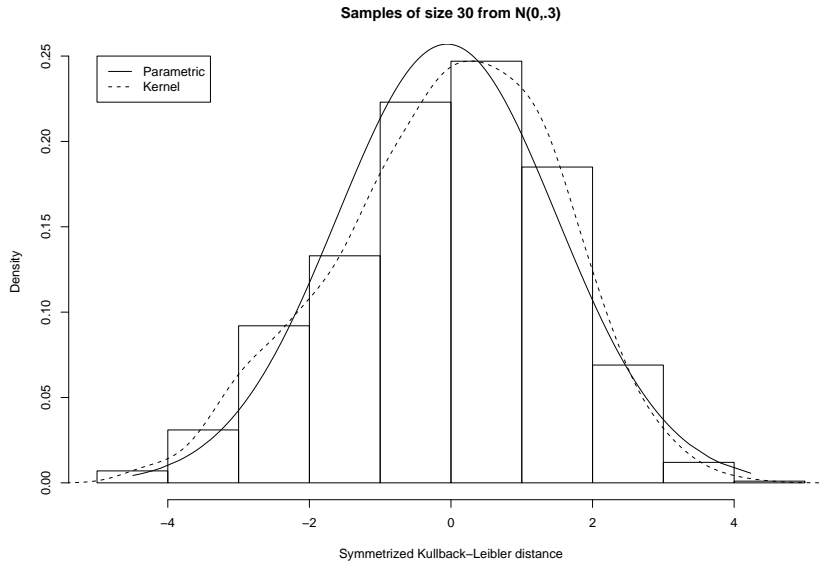


Figure 1 *Parametric and kernel estimate for the distribution of SKL*

Table 1 *Power function for SKL and K-S for mixture of normal distributions*

μ	-.5	0	.5	.7	.9	1.0	1.1	1.2	1.5
SKL	.045	.099	.161	.256	.384	.541	.663	.782	.991
K-S	.045	.098	.124	.233	.356	.444	.500	.616	.885

5 Numerical example

In this section we use data from two sources. One of them is the data from Ramus et al. (1999). The other consists of 20 sentences from Portuguese read by two speakers of Modern European Portuguese and Brazilian Portuguese (EP and BP respectively). These sentences were designed by Sonia Frota and Charlotte Galves to study several characteristics of Portuguese, not only consonantal and vocalic intervals, but also stressed syllables, secondary stressed syllables among others. These sentences were recorded at 16 kHz and 11kHz and then segmented by hand by two persons. They used both audio and visual clues to identify consonantal and vocalic intervals and used Multi Speech 3700 software to analyze the acoustic

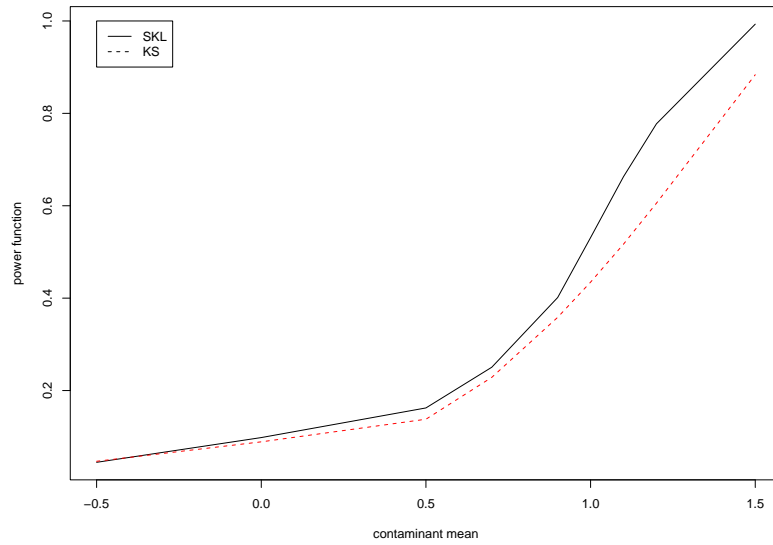


Figure 2 *Power function for SKL and K-S for mixture of normal distributions*

signal. The same procedure was used by Ramus et al. (1999) to record and segment the other acoustic data as well.

For each sentence the duration of the consonantal intervals were computed, call this variable C . This variable is important in view of the work of Ramus et al. (1999) which could cluster 8 languages into into three groups which correspond precisely to the intuitive notion of rhythmic classes: English, Polish and Dutch represent the accentual class, French, Spanish, Catalan and Italian represent the syllabic class and Japanese represents the moraic class. The same variable was used by Duarte, Galves, Garcia and Maronna (2001) using a parametric approach adjusting a gamma model to fit the data from all languages. Maximum likelihood ratio tests seems to confirm Ramus et al. classification and placed European Portuguese among the accentual languages and Brazilian Portuguese among the syllabic ones. Using SKL and K-S to compare the distribution of C for some of the languages we obtained somehow different results. First of all, we cannot distinguish between Brazilian and European Portuguese (p-values: SKL=.69 and K-S=.66). Also, at a 5% significance level there is evidence of difference between Brazilian Portuguese and Catalan (p-values: SKL=.02 and K-S=.02). Figure 3 presents density estimates by kernel and by SKL (Dias, 2000) suggesting that Catalan is bimodal, maybe a mixture of two gammas and this causes the tests to reject the equality of the distributions.

As conjectured there is significant evidence for difference between English and

Japanese (p-values: SKL=.01 and K-S $< 10^{-3}$) and no evidence of difference between English and Dutch (p-values: SKL=.59 and K-S=.18) and English and European Portuguese (p-values: SKL=.15 and K-S=.05).

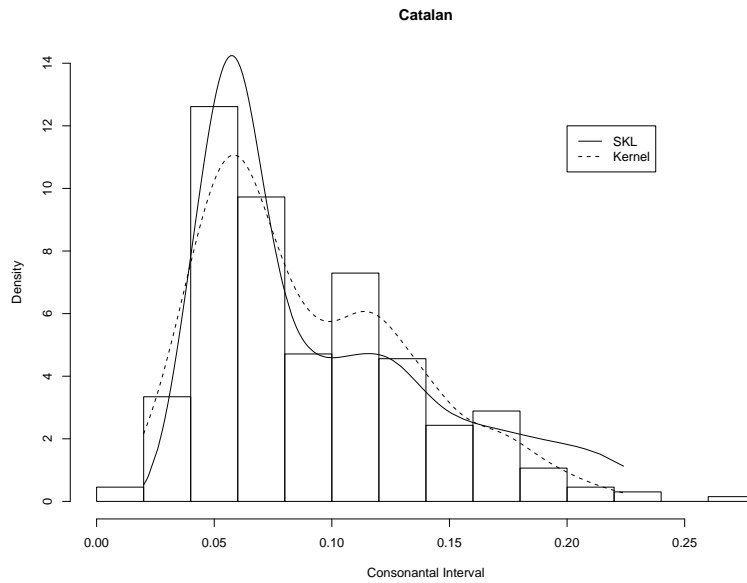


Figure 3 *Density estimates for consonantal intervals of Catalan*

Acknowledgments

We thank Ricardo Molina for recording the data, Janaisa Viscardi and Flaviane Fernandes for segmenting the data, Frank Ramus, Marina Nespór and Jacques Meller for making their data available. We also thank Denise Duarte, Antonio Galves and Charlotte Galves for many fruitful discussions. The presentation of this paper was improved by valuable remarks from the anonymous referee. This work is partially funded by FAPESP grant 02/01554-5 and CNPq grants 301054/1993-2, 300644/1994-9 and 475763/2003-3.

(Received September, 2003. Accepted June, 2004.)

References

Abercrombie, D. (1967). *Elements of general phonetics*. Chicago: Aldine.

- Dias, R. (1994). *Density estimation via H-splines*. University of Wisconsin-Madison, Ph.D. dissertation.
- Dias, R. (1998). Density estimation via hybrid splines. *Journal of Statistical Computation and Simulation*, **60**, 277–294.
- Dias, R. (2000). A note on density estimation using a proxy of the Kullback-Leibler distance. *Brazilian Journal of Probability and Statistics*, **13**, 181–192.
- Dias, R. and Gamerman, D. (2002). A Bayesian approach to Hybrid splines non-parametric regression. *Journal of Statistical Computation and Simulation*, **72**, 285–297.
- Duarte, D., Galves, A., Garcia, N. L. and Maronna, R. (2001). The statistical analysis of acoustic correlates of speech rhythm. *Workshop on Rhythmic Patterns, Parameter Setting and Language Change*, ZiF, University of Bielefeld. <http://www.physik.uni-bielefeld.de/complexity/duarte.pdf>. Paper presented.
- Fan, Y. (1998). Goodness-of-fit tests based on kernel density estimators with fixed smoothing parameters. *Econometric Theory*, **14**, 604–621.
- Ferguson, T. S. (1996). A course in large sample theory. *Texts in Statistical Science Series*. London: Chapman & Hall.
- Hall, P., Nussbaum, M. and Stern, S. E. (1997). On the estimation of a support curve of indeterminate sharpness. *J. Multivariate Anal.*, **62**, 204–232.
- Kooperberg, C. and Stone, C. J. (1991). A study of logspline density estimation. *Computational Statistics and Data Analysis*, **12**, 327–347.
- Li, Qi (1996). Nonparametric testing of closeness between two unknown distribution functions. *Econometric Reviews*, **15**, 261–274.
- Mason, D. M. and Schuenemeyer, J. H. (1983). A modified Kolmogorov-Smirnov test sensitive to tail alternatives. *Ann. Statist.*, **11**, 933–946.
- Mason, D. M. and Schuenemeyer, J. H. (1992). Correction: “A modified Kolmogorov-Smirnov test sensitive to tail alternatives” [*Ann. Statist.*, **11**, (1983), 933–946; MR 85c:62113]. *Ann. Statist.*, **20**, 620–621.
- Pagan, A. and Ullah, A. (1999). *Nonparametric Econometrics*. Cambridge: Cambridge University Press.
- Pike, K. L. (1945). *The Intonation of American English*. Ann Arbor: University of Michigan Press.
- Ramus, F., Nespors, M. and Mehler, J. (1999). Correlates of linguistic rhythm in speech. *Cognition*, **73**, 265–292.

Silverman, B. W. (1986). *Density Estimation for Statistics and Data Analysis*. 184p. London: Chapman and Hall.

Vidakovic, B. (1999). Statistical modelling by wavelets. *Wiley Series in Probability and Statistics: Applied Probability and Statistics*. A Wiley-Interscience Publication. New York: John Wiley and Sons Inc.

Ronaldo Dias and **Nancy L. Garcia**

Departamento de Estatística,
Universidade Estadual de Campinas
São Paulo, Brasil
E-mails: dias@ime.unicamp.br and
nancy@ime.unicamp.br