

# Estimating Mixing Densities in Exponential Family Models for Discrete Variables

WEI-LIEM LOH

*Purdue University and National University of Singapore*

CUN-HUI ZHANG

*Rutgers University*

**ABSTRACT.** This paper is concerned with estimating a mixing density  $g$  using a random sample from the mixture distribution  $f(x) = \int f(x | \theta)g(\theta) d\theta$  where  $f(\cdot | \theta)$  is a known discrete exponential family of density functions. Recently two techniques for estimating  $g$  have been proposed. The first uses Fourier analysis and the method of kernels and the second uses orthogonal polynomials. It is known that the first technique is capable of yielding estimators that achieve (or almost achieve) the minimax convergence rate. We show that this is true for the technique based on orthogonal polynomials as well. The practical implementation of these estimators is also addressed. Computer experiments indicate that the kernel estimators give somewhat disappointing finite sample results. However, the orthogonal polynomial estimators appear to do much better. To improve on the finite sample performance of the orthogonal polynomial estimators, a way of estimating the optimal truncation parameter is proposed. The resultant estimators retain the convergence rates of the previous estimators and a Monte Carlo finite sample study reveals that they perform well relative to the ones based on the optimal truncation parameter.

*Key words:* discrete exponential family, mixing density, Monte Carlo simulation, orthogonal polynomials, rate of convergence

## 1. Introduction

Let  $X_1, \dots, X_n$  be independent observations from a mixture distribution with probability law

$$f(x; g) = \int_0^{\theta^*} f(x | \theta)g(\theta) d\theta, \quad (1)$$

where  $g$  is a mixing probability density function on  $(0, \theta^*)$  and  $f(\cdot | \theta)$  is a known parametric family of probability density functions with respect to a  $\sigma$ -finite measure  $\nu$ . In particular we assume that

$$f(x | \theta) = C(\theta)q(x)\theta^x, \quad \forall x = 0, 1, 2, \dots, \quad (2)$$

where  $0 < \theta < \theta^* \leq \infty$ ,  $q(x) > 0$  whenever  $x = 0, 1, 2, \dots$  and  $\nu$  is the counting measure on the set of non-negative integers. In this paper we are concerned with the estimation of  $g$  using the random sample  $X_1, \dots, X_n$ .

Over the last few years, there has been a great deal of interest in the above problem and over related mixture problems. Important advances have been made on the deconvolution problem by Devroye & Wise (1979), Carroll & Hall (1988), Zhang (1990), Fan (1991) and many others using Fourier techniques. In particular kernel estimators have been obtained which achieve the minimax convergence rate.

In the context of mixtures of discrete exponential families, Tucker (1963) considered the estimation of the mixing distribution of a Poisson mixture via the method of moments and Simar (1976) approached the same problem using maximum likelihood. Rolph (1968), Meeden (1972) and Datta (1991) used Bayesian methods to construct consistent estimators for the mixing distribution.

Quite recently, two techniques for the estimation of the mixing density  $g$ , as given in (1), have been proposed. The first was proposed by Zhang (1995) which uses Fourier analysis and the method of kernels. The second was proposed by Walter & Hamedani (1989, 1991) which uses orthogonal polynomials. It has been shown by Zhang (1995) and Loh & Zhang (1996) that the first technique is capable of yielding estimators that achieve (or almost achieve) the minimax convergence rate with respect to local and integrated mean squared error over various smoothness classes of mixing density functions.

*Remark 1.* During the revision of this paper, we became aware of Hengartner (1995) who showed that for a Poisson mixture with  $\theta^* < \infty$ , orthogonal polynomial mixing density estimators can attain the minimax convergence rate with respect to integrated mean squared error over the Sobolev space of mixing density functions with square integrable  $r$ th derivatives.

The rest of this paper is organized as follows. We shall first very briefly review the kernel mixing density estimators and their properties in section 2. In section 3 we shall show that the technique based on orthogonal polynomials is also capable of yielding mixing density estimators that achieve (or almost achieve) the minimax convergence rate with respect to integrated weighted mean squared error over various non-parametric classes of mixing density functions. However, even with this property the minimax convergence rates of these estimators are logarithmic (not polynomial). This leaves us with the important question as to how well can these estimators actually perform in practice.

Section 4 addresses the issue of the finite sample performances as well as the practical implementation of these estimators. Computer experiments indicate that the kernel mixing density estimators (for the particular kernel used here) give somewhat disappointing finite sample results. On the other hand, the orthogonal polynomial mixing density estimators appear to do much better. To improve upon the finite sample performance of the orthogonal polynomial mixing density estimators further, a way of estimating the optimal truncation parameter is proposed in section 5. The resultant estimators retain the convergence rates of the previous estimators and a Monte Carlo finite sample study reveals that they perform well relative to the ones based on the optimal truncation parameter.

All proofs in this paper have been deferred to the appendix. Finally we shall denote by  $P = P_g$  and  $E = E_g$  the probability and expectation corresponding to  $g$  respectively, by  $h^{(j)}$  the  $j$ th derivative (if it exists) of any function  $h$  with  $h^{(0)} = h$ , and the weighted  $L^p$ -norm of any measurable function  $h$  by  $\|h\|_{w,p} = (\int |h(y)|^p w(y) dy)^{1/p}$ ,  $\forall 1 \leq p < \infty$ . If  $w(y) \equiv 1$ , we denote  $\|\cdot\|_{w,p}$  by  $\|\cdot\|_p$ .

## 2. Kernel mixing density estimators

This section treats the case  $\theta^* < \infty$  and, for completeness, gives a brief review of the kernel mixing density estimators that we are concerned with here. We refer the reader to Loh & Zhang (1996) for the proofs and a more detailed discussion for these estimators.

Let  $k: \mathcal{R} \rightarrow \mathcal{R}$  be a symmetric function satisfying

$$\int_{-\infty}^{\infty} k(y) dy = 1, \quad k^*(t) = 0, \quad \forall |t| > 1,$$

$$\int_{-\infty}^{\infty} y^j k(y) dy = 0, \quad \forall j < \alpha_0, \quad (3)$$

and

$$\int_{-\infty}^{\infty} |y^{\alpha_0} k(y)| dy < \infty, \tag{4}$$

for some positive number  $\alpha_0$ , where  $k^*$  denotes the Fourier transform of  $k$ , that is

$$k^*(t) = \int_{-\infty}^{\infty} \exp(ity) k(y) dy.$$

Define

$$K_n(x, \theta) = \frac{I\{0 \leq x \leq d_n\}}{2\pi q(x)x!} \int_{-c_n}^{c_n} \Re\{(it)^x \exp(-it\theta)\} k^*(t/c_n) dt, \tag{5}$$

where  $c_n$  and  $d_n$  are positive constants tending to  $\infty$ ,  $I\{\cdot\}$  denotes the indicator function and  $\Re(z)$  is the real part of the complex number  $z$ . Observing that

$$E_g K_n(X_1, \theta) - C(\theta)g(\theta) \rightarrow 0, \quad \forall -\infty < \theta < \infty, \tag{6}$$

as  $(c_n, d_n) \rightarrow (\infty, \infty)$  along a suitable path, we may estimate  $g(\theta)$  by the kernel mixing density estimator

$$\hat{g}_{K,n}(\theta) = n^{-1} \sum_{j=1}^n \{K_n(X_j, \theta)/C(\theta)\} I\{0 < \theta < a_n\}, \quad \forall 0 < \theta < \theta^*, \tag{7}$$

where  $a_n, c_n$ , and  $d_n$  are constants satisfying

$$c_n + \max_{1 \leq x \leq d_n} \log(1/q(x)) = \beta_0 \log n, \quad c_n = (\theta^* e)^{-1} (d_n - \beta_1 \log c_n), \tag{8}$$

and

$$a_n = \begin{cases} \theta^* & \text{if } C(\theta^*) > 0, \\ \theta^* - a^*/c_n & \text{if } C(\theta^*) = 0, \end{cases} \tag{9}$$

with absolute (independent of  $n$ ) constants  $0 < \beta_0 < 1/2$ ,  $\beta_1 > 0$ , and  $0 < a^* < \infty$ . The performance of these estimators is investigated with respect to the following smoothness classes of mixing density functions. Let  $w$  be a measurable function on  $(0, \theta^*)$  with  $\|w\|_1$  finite. For  $\alpha > 0$  we define  $\mathcal{G}_{\alpha, \theta^*}(w, M)$  to be the set of all probability density functions  $g$  on  $(0, \theta^*)$  such that

$$\|g^{(\alpha')}(\cdot) - g^{(\alpha')}(\cdot + \delta)\|_{w,2} < M|\delta|^{\alpha'}, \quad \forall \delta, \tag{10}$$

where  $\alpha'$  is the integer with  $0 < \alpha'' = \alpha - \alpha' \leq 1$ , and  $M$  is a constant such that  $\mathcal{G}_{\alpha, \theta^*}(w, M)$  is non-empty.

We further assume that there exist constants  $\gamma \geq 0$ ,  $C_1^*$ ,  $C_2^*$ , and  $C_3^*$  such that

$$\sup_{0 < \theta < \theta^*} (\theta^* - \theta)^\gamma / C(\theta) < C_1^*, \tag{11}$$

$$\sup_{0 < \theta < \theta^*} (\theta^* - \theta)^\gamma |C^{(j)}(\theta)| / \{C(\theta)j!\} < C_2^*, \quad \forall 0 \leq j \leq \rho', \tag{12}$$

and

$$|C^{(\rho')}(\theta + \delta) - C^{(\rho')}(\theta)| < C_3^* \delta^{\rho'}, \quad 0 < \theta < \theta + \delta < \theta^*, \tag{13}$$

where  $\rho'$  is a non-negative integer with  $0 < \rho'' = \rho - \rho' \leq 1$ .

Theorem 1 below shows that the kernel mixing density estimators  $\hat{g}_{K,n}$  achieve (or almost achieve) the minimax convergence rate with respect to  $\mathcal{G}_{\alpha, \theta^*}(w, M)$  under reasonably mild conditions.

**Theorem 1**

Suppose  $\alpha > 0$  and that (11)–(13) hold with  $\gamma \geq 0$  and  $\rho = \alpha + \gamma$ . Let  $\hat{g}_{K,n}$  be given by (7) with the kernel  $K_n(x, \theta)$  in (5) such that  $\alpha_0 \geq \alpha + \gamma$  in (4). Let (8) and (9) hold with  $\beta_1 \leq \alpha + \gamma$ . Then if

$$q(x)\gamma_0\gamma_1^x(x!)^\beta \geq 1, \quad \forall x \geq 0,$$

for some constants,  $\gamma_0$ ,  $\gamma_1$ , and  $\beta$ , we have

$$\sup_{g \in \mathcal{G}_{\alpha, \theta^*}(w, M)} E_g \|\hat{g}_{K,n} - g\|_{w,2} = \begin{cases} O(1)(1/\log n)^\alpha & \text{if } \beta = 0, \\ O(1)(\log \log n / \log n)^\alpha & \text{if } 0 < \beta < \infty. \end{cases}$$

Furthermore, if  $g_0$  is an interior point of  $\mathcal{G}_{\alpha, \theta^*}(1, M)$ , then

$$\liminf_{n \rightarrow \infty} (\log n)^\alpha \inf_{\hat{g}_n} \sup \{E_g \|\hat{g}_n - g\|_2 : g \in \mathcal{G}_{\alpha, \theta^*}(1, M), \|g - g_0\|_2 \leq M_1(\log n)^{-\alpha}\} > 0,$$

where the infimum runs over all possible estimators  $\hat{g}_n$  based on  $X_1, \dots, X_n$ ,  $M_1$  is a positive constant and  $\mathcal{G}_{\alpha, \theta^*}(1, M)$  is given by (10) with  $w(\theta) = I\{0 < \theta < \theta^*\}$ .

**3. Orthogonal polynomial mixing density estimators**

In this section we introduce the class of orthogonal polynomial mixing density estimators that we are concerned with and also establish upper and lower bounds for their convergence rates with respect to various non-parametric classes of mixing density functions. Let  $C: (0, \theta^*) \rightarrow R^+$  be as in (2) and  $w: (0, \theta^*) \rightarrow R^+$  be a measurable function such that  $\|C^2/w\|_1 < \infty$ . Let  $\{p_{w_0,j}\}_{j=0}^\infty$  be a sequence of orthogonal polynomials on  $(0, \theta^*)$  with weight function

$$w_0(\theta) = C^2(\theta)/w(\theta). \quad (14)$$

In particular, we assume that these polynomials are normalized so that

$$p_{w_0,j}(\theta) = \sum_{x=0}^j k_{w_0,j,x} \theta^x, \quad (15)$$

with  $k_{w_0,j,j} > 0$  for all  $j \geq 0$ , and  $\int_0^{\theta^*} p_{w_0,i}(\theta)p_{w_0,j}(\theta)w_0(\theta) d\theta = \delta_{ij}$ , where  $\delta_{ij}$  denotes the Kronecker delta. We further assume that  $\{p_{w_0,j}\}_{j=0}^\infty$  is complete with respect to  $\|\cdot\|_{w_0,2}$ . Note that this is always true if  $\theta^* < \infty$  [see for example Szegő (1975) p. 40]. Next define

$$\lambda_{w_0,j}(x) = \begin{cases} k_{w_0,j,x}/q(x) & \text{if } 0 \leq x \leq j, \\ 0 & \text{otherwise.} \end{cases}$$

We write

$$h(\theta) = w(\theta)g(\theta)/C(\theta), \quad \forall 0 < \theta < \theta^*, \quad (16)$$

and assume that the mixing density  $g$  satisfies  $\|g\|_{w,2} = \|h\|_{w_0,2} < \infty$ . Then  $h$  has the formal orthogonal polynomial series expansion  $h(\theta) \sim \sum_{j=0}^\infty h_{w_0,j} p_{w_0,j}(\theta)$ , where

$$h_{w_0,j} = \int_0^{\theta^*} h(\theta)p_{w_0,j}(\theta)w_0(\theta) d\theta, \quad \forall j = 0, 1, 2, \dots \quad (17)$$

Observing that

$$E_g \lambda_{w_0,j}(X_1) = \sum_{x=0}^\infty f(x; g) \lambda_{w_0,j}(x) = h_{w_0,j}, \quad \forall j = 0, 1, 2, \dots,$$

we estimate  $h_{w_0,j}$  by  $\hat{h}_{w_0,j} = n^{-1} \sum_{i=1}^n \lambda_{w_0,j}(X_i)$  and  $g(\theta)$  by the orthogonal polynomial mixing density estimator

$$\hat{g}_{OP,n}(\theta) = [C(\theta)/w(\theta)] \sum_{j=0}^{m_n} \hat{h}_{w_0,j} p_{w_0,j}(\theta), \quad \forall 0 < \theta < \theta^*, \tag{18}$$

where  $m_n$  is a positive constant (truncation parameter) which tends to  $\infty$  as  $n \rightarrow \infty$ . The following proposition gives an upper bound on the convergence rate of  $\hat{g}_{OP,n}$ .

**Proposition 1**

Suppose  $\|C^2/w\|_1 < \infty$  and  $\|g\|_{w,2} < \infty$ . Let  $\hat{g}_{OP,n}$  be as in (18). Then

$$E_g \|\hat{g}_{OP,n} - g\|_{w,2}^2 \leq n^{-1} \sum_{j=0}^{m_n} \max_{0 \leq x \leq j} [k_{w_0,j,x}/q(x)]^2 + \sum_{j=m_n+1}^{\infty} h_{w_0,j}^2,$$

with  $k_{w_0,j,x}$  and  $h_{w_0,j}$  as in (15) and (17) respectively.

*Remark 2.* The motivation for (18) originates from Walter & Hamedani (1989) who proposed a similar class of estimators. They also obtained a result analogous to proposition 1.

We now study the performance of the estimators  $\hat{g}_{OP,n}$  with respect to the following non-parametric classes of mixing density functions. For positive constants  $\alpha, M$  and  $m = 1, 2, \dots$ , we define  $\mathcal{G}(\alpha, m, M, w_0)$  to be the set of all probability density functions  $g$  on  $(0, \theta^*)$  such that  $\|g\|_{w,2} < \infty$  and  $\sum_{j=m}^{\infty} j^{2\alpha} h_{w_0,j}^2 < M$  with  $h_{w_0,j}$  as in (17). We note that this class implicitly depends on the discrete exponential family of interest, in particular on  $C(\theta)$ . This ellipsoidal class is chosen mainly for reasons of mathematical tractability. However, ellipsoid conditions can amount to the imposition of smoothness and integrability requirements, see for example Johnstone & Silverman (1990) p. 258. In our case, we have the following characterization.

**Proposition 2**

Let  $m \geq 1$  and  $\{p_{w_0,j}\}_{j=0}^{\infty}$  be as in (15). Suppose there exist constants  $v_{j,m}, j \geq m$  and another sequence of (normalized) complete orthogonal polynomials  $\{p_{w_1,j}\}_{j=0}^{\infty}$  with weight function  $w_1$  such that

$$[p_{w_1,j}(\theta)w_1(\theta)]^{(m)} = (-1)^m v_{j+m,m} p_{w_0,j+m}(\theta)w_0(\theta), \quad \forall j \geq 0, \tag{19}$$

and

$$\alpha_1 < \inf_{j \geq m} |v_{j,m}|/j^\alpha \leq \sup_{j \geq m} |v_{j,m}|/j^\alpha < \alpha_2, \tag{20}$$

where  $\alpha, \alpha_1$  and  $\alpha_2$  are positive constants. If  $\|h^{(m)}\|_{w_1,2} < \infty$  and

$$0 = \lim_{\theta \rightarrow 0^+} h^{(m-i)}(\theta)[p_{w_1,j}(\theta)w_1(\theta)]^{(i-1)} = \lim_{\theta \rightarrow \theta^*-} h^{(m-i)}(\theta)[p_{w_1,j}(\theta)w_1(\theta)]^{(i-1)} \tag{21}$$

for  $0 < i \leq m$  and  $j \geq 0$ , then

$$\alpha_1 \left( \sum_{j=m}^{\infty} j^{2\alpha} h_{w_0,j}^2 \right)^{1/2} \leq \|h^{(m)}\|_{w_1,2} \leq \alpha_2 \left( \sum_{j=m}^{\infty} j^{2\alpha} h_{w_0,j}^2 \right)^{1/2}, \tag{22}$$

where  $h_{w_0,j}$  is defined as in (17).

**Remark 3.** It is shown in the appendix that (19) and (20) satisfied by (i) the classical orthogonal polynomials of Laguerre with  $\alpha = m/2$  and  $w_0(\theta) = \theta^\beta e^{-\theta}$ ,  $\theta > 0$  and (ii) the classical orthogonal polynomials of Jacobi with  $\alpha = m$  and  $w_0(\theta) = \theta^{\beta_1}(\theta^* - \theta)^{\beta_2}$ , where  $\beta > -1$  and  $\beta_j > -1, j = 1, 2$ . In general,  $\alpha$  and  $m$  become functions of each other under (19) and (20), although they are two independent parameters in the definition of  $\mathcal{G}(\alpha, m, M, w_0)$ .

For the rest of this section, we shall assume that  $M$  is sufficiently large so that  $\mathcal{G}(\alpha, m, M, w_0)$  is non-empty. The next two theorems and their corollaries establish upper bounds on the convergence rate of  $\hat{g}_{OP, n}$  over the class of mixing densities  $\mathcal{G}(\alpha, m, M, w_0)$ .

### Theorem 2

Suppose  $\|C^2/w\|_1 < \infty$ . Let  $\hat{g}_{OP, n}$  be as in (18) and

$$\max_{0 \leq x \leq j \leq m_n} \log(|k_{w_0, j, x}|/q(x)) \leq \beta_0 \log n, \quad (23)$$

for some constant  $0 < \beta_0 < 1/2$  where  $w_0$  is as in (14). Then

$$\sup \{E_g \|\hat{g}_{OP, n} - g\|_{w, 2} : g \in \mathcal{G}(\alpha, m, M, w_0)\} = O(1)(m_n^{-\alpha} + m_n^{1/2} n^{(2\beta_0 - 1)/2}).$$

### Corollary 1

Suppose  $\theta^* = \infty$ ,  $w(\theta) = \theta^{-\beta} C^2(\theta) e^\theta$  and  $w_0(\theta) = \theta^\beta e^{-\theta}$  with  $\beta > -1$ . Let  $\{p_{w_0, j}\}_{j=0}^\infty$  be the sequence of (normalized) Laguerre polynomials on  $(0, \infty)$  with weight function  $w_0$ ,  $\hat{g}_{OP, n}$  as in (18) and

$$q(x)\gamma_0\gamma_1^\gamma(x!) > 1, \quad \forall x \geq 0,$$

for constants  $\gamma_0$  and  $\gamma_1$ . Then by choosing  $m_n = \delta \log n$  with  $0 < \delta \leq \beta_0/\log(2\gamma_1)$  and  $0 < \beta_0 < 1/2$ , we have

$$\sup \{E_g \|\hat{g}_{OP, n} - g\|_{w, 2} : g \in \mathcal{G}(m/2, m, M, w_0)\} = O(1)(1/\log n)^{m/2}.$$

Theorem 3 is a specialization of theorem 2 which proves to be useful when  $\theta^* < \infty$ .

### Theorem 3

Let  $\hat{g}_{OP, n}$  be as in (18) and that for some constant  $\zeta > 1$ ,

$$\max_{0 \leq x \leq j} k_{w_0, j, x}^2 < \zeta^{2j}, \quad \forall j \geq 0, \quad (24)$$

where  $w_0$  is as in (14). Suppose further that

$$\max_{0 \leq x \leq m_n} \log(1/q(x)) + m_n \log \zeta \leq \beta_0 \log n, \quad (25)$$

with constant  $0 < \beta_0 < 1/2$ . Then

$$\sup \{E_g \|\hat{g}_{OP, n} - g\|_{w, 2} : g \in \mathcal{G}(\alpha, m, M, w_0)\} = O(m_n^{-\alpha}).$$

### Corollary 2

Let  $\hat{g}_{OP, n}$  be as in (18) and that (24) holds for some constant  $\zeta > 1$ . Suppose

$$q(x)\gamma_0\gamma_1^\gamma(x!)^\gamma > 1, \quad \forall x \geq 0, \quad (26)$$

for constants  $\gamma_0, \gamma_1 \geq 1$  and  $\gamma$ . Then

(a) if  $\gamma = 0$ , by choosing  $m_n = \delta \log n$  with  $0 < \delta \leq \beta_0/\log(\gamma_1\zeta)$  and  $0 < \beta_0 < 1/2$ , we have

$$\sup \{E_g \|\hat{g}_{OP, n} - g\|_{w, 2} : g \in \mathcal{G}(\alpha, m, M, w_0)\} = O(1)(1/\log n)^\alpha,$$

(b) if  $0 < \gamma < \infty$ , by choosing  $m_n = \delta \log n / \log \log n$  with  $0 < \delta < \beta_0 / \gamma$  and  $0 < \beta_0 < 1/2$ , we have

$$\sup \{E_g \|\hat{g}_{OP, n} - g\|_{w, 2} : g \in \mathcal{G}(\alpha, m, M, w_0)\} = O(1)(\log \log n / \log n)^\alpha.$$

*Remark 4.* The negative binomial and Poisson mixtures satisfy (26) with  $\gamma = 0$  and 1 respectively.

*Remark 5.* The classical orthogonal polynomials of Jacobi satisfy (24).

The next theorem complements the above results by establishing lower bounds on the local minimax convergence rate over the class of mixing densities  $\mathcal{G}(\alpha, m, M, w_0)$  under the condition that (19) and (20) hold.

#### Theorem 4

Let  $w: (0, \theta^*) \rightarrow R^+$  be a measurable function such that  $\|w\|_1 < \infty$  and  $\|w_0\|_1 < \infty$  with  $w_0$  as in (14) and  $\{p_{w_0, j}\}_{j=0}^\infty$  be a sequence of (normalized) orthogonal polynomials with weight function  $w_0$  such that (19) and (20) are satisfied. Suppose there exists an open interval where  $w$  is strictly positive and  $m$  times continuously differentiable. Then for sufficiently large  $M$ , we have

$$\liminf_{n \rightarrow \infty} (\log n)^m \inf_{\hat{g}_n} \sup \{E_g \|\hat{g}_n - g\|_{w, 2} : g \in \mathcal{G}(\alpha, m, M, w_0), \|g - g_0\|_{w, 2} \leq M(\log n)^{-m}\} > 0$$

for each interior point  $g_0$  of  $\mathcal{G}(\alpha, m, M, w_0)$ , where the infimum runs over all possible estimators  $\hat{g}_n$  based on  $X_1, \dots, X_n$ .

Following Hengartner (1995), the next theorem gives conditions for the sharpening of the lower bounds of theorem 4 when  $\theta^* < \infty$ .

#### Theorem 5

Let  $\theta^* < \infty$ ,  $w: (0, \theta^*) \rightarrow R^+$  be a measurable function such that  $\|w\|_1 < \infty$  and  $\|w_0\|_1 < \infty$  with  $w_0$  as in (14) and  $\{p_{w_0, j}\}_{j=0}^\infty$  be a sequence of (normalized) orthogonal polynomials with weight function  $w_0$  such that (19) and (20) are satisfied. Suppose

$$q(x)(x!)^\tau < \tau_0 \tau_1^x, \quad \forall x \geq 0, \quad (27)$$

for strictly positive constants  $\tau_0$ ,  $\tau_1$  and  $\tau$  and that there exists an open interval where  $w$  is strictly positive and  $m$  times continuously differentiable. Then for sufficiently large  $M$ , we have

$$\liminf_{n \rightarrow \infty} (\log n / \log \log n)^m \inf_{\hat{g}_n} \sup \{E_g \|\hat{g}_n - g\|_{w, 2} : g \in \mathcal{G}(\alpha, m, M, w_0)\} > 0,$$

where the infimum runs over all possible estimators  $\hat{g}_n$  based on  $X_1, \dots, X_n$ .

We close this section with the following consequence of corollary 2, remarks 3 and 5 and theorems 4 and 5. Suppose  $\theta^* < \infty$ . Let

$$w(\theta) = C^2(\theta)\theta^{-\beta_1}(\theta^* - \theta)\theta^{-\beta_2}, \quad \forall 0 < \theta < \theta^*,$$

with  $\beta_j > -1$ ,  $j = 1, 2$ , so that (19) and (20) hold for  $w_0(\theta) = \theta^{\beta_1}(\theta^* - \theta)^{\beta_2}$  and the Jacobi polynomials  $\{p_{w_0, j}\}_{j=0}^\infty$ .

If  $q(x)\gamma_0\gamma_1^x > 1 \forall x \geq 0$  for constants  $\gamma_0$  and  $\gamma_1$ , then the minimax convergence rate with respect to  $\|\cdot\|_{w, 2}$  loss is  $(1/\log n)^m$  for mixing densities  $g$  in the class  $\mathcal{G}(m, m, M, w_0)$  where  $w_0$  is as in (14).

On the other hand, if  $q(x)\gamma_0\gamma_1^x(x!)^\gamma > 1 \forall x \geq 0$  and  $q(x)(x!)^\tau < \tau_0\tau_1^x \forall x \geq 0$  for constants  $\gamma_0, \gamma_1, \tau_0, \tau_1$  and  $0 < \gamma, \tau < \infty$ , then the minimax convergence rate with respect to  $\|\cdot\|_{w,2}$  loss is  $(\log \log n / \log n)^m$  for mixing densities  $g$  in the class  $\mathcal{G}(m, m, M, w_0)$ . We note that both minimax rates are attained by the orthogonal polynomial mixing density estimators  $\hat{g}_{OP,n}$  of corollary 2.

#### 4. Finite sample performance

A key consequence of the results of sections 2 and 3 is that both the kernel and orthogonal polynomial mixing density estimators, that is  $\hat{g}_{K,n}$  and  $\hat{g}_{OP,n}$  respectively, are capable of achieving (or almost achieving) the minimax rate of convergence. However, even with this property the minimax convergence rate of these estimators is logarithmic (not polynomial). This leads us to the following problem: typically how large must a sample be in order that the desired asymptotics of these estimators (as described in the previous two sections) can take effect.

##### 4.1. Kernel mixing density estimators

In order to gauge typically how well the kernel mixing density estimators perform in practice, we focus on the problem of estimating the mixing density  $g$  of a negative binomial mixture with  $\theta^* = 1$  and  $C(\theta) = 1 - \theta$  with respect to integrated squared error, that is  $\|\hat{g}_n - g\|_2^2$ . To construct the kernel mixing density estimator  $\hat{g}_{K,n}$ , we take

$$k(y) = \frac{6}{\pi} \frac{2}{y} \sin\left(\frac{y}{4}\right)^4, \quad \forall -\infty < y < \infty.$$

Our motivation for such a choice of  $k$  is its relative simplicity and that (3) and (4) hold with  $\alpha_0 = 2$ . We observe from (6) and (7) that an upper bound on the finite sample performance of  $\hat{g}_{K,n}$  can be obtained by investigating how close

$$\|E_g[K_n(X_1, \theta)/C(\theta)]I\{0 < \theta < a_n\} - g(\theta)\|_2^2 \quad (28)$$

is to 0. In this case we take  $g(\theta) = I\{0 < \theta < 1\}$  and use

$$ERR_n = (1/10) \sum_{i=1}^{10} \{E_g[K_n(X_1, 0.1i - 0.05)/C(0.1i - 0.05)] - 1\}^2$$

as an approximation to (28).

*Remark 6.* The reason for such a choice of  $g$  is that we feel that the uniform distribution is arguably one of the distributions that any reasonable estimation procedure should be able to estimate adequately well.

Computations show that in order to have  $ERR_n \approx 0.1$ , we need  $c_n \approx 17$ . Since  $c_n \leq (1/2) \log n$ , this implies that the sample size  $n$  must be astronomically large and is quite impossible to obtain in practice.

This presents a disappointing setback for the practical implementation of  $\hat{g}_{K,n}$ . However, it should be noted that this can be due to a possibly inappropriate choice of the kernel  $k$  and that it does not eliminate the possibility that there exist other kernels which give dramatically better results.



4.2. Orthogonal polynomial mixing density estimators

We observe that the integrated mean squared error of the orthogonal polynomial mixing density estimators has a simple closed form expression. In particular, we observe as in (32) that

$$E_g \int_0^{\theta^*} [\hat{g}_{OP,n}(\theta) - g(\theta)]^2 w(\theta) d\theta = \int_0^{\theta^*} g^2(\theta) w(\theta) d\theta + n^{-1} \sum_{j=0}^{m_n} \{E_g \lambda_{w_0,j}^2(X_1) - (n+1)[E_g \lambda_{w_0,j}(X_1)]^2\}. \tag{29}$$

The right hand side of (29) enables us to compute the integrated mean squared error of  $\hat{g}_{OP,n}$  in any given situation. We illustrate this below with two examples.

*Example 1.* This example deals with the problem of estimating a mixing density  $g$  of a negative binomial mixture with  $\theta^* = 1$  and  $C(\theta) = 1 - \theta$  using integrated squared error loss. In this case the orthogonal polynomial mixing density estimators are given as in (18) where  $\{p_{w_0,j}\}_{j=0}^{\infty}$  corresponds to the Jacobi polynomials with weight function  $w_0(\theta) = (1 - \theta)^2$ ,  $\forall 0 < \theta < 1$ .

Tables 1, 2 and 3 give the integrated mean squared error of  $\hat{g}_{OP,n}$  for sample sizes  $n = 1000, 10\ 000$  and  $100\ 000$  as well as for truncation parameters  $0 \leq m_n \leq 4$ .

*Example 2.* This example deals with the estimation of the mixing density  $g$  of a Poisson mixture with  $\theta^* = \infty$  using integrated weighted squared error loss  $\|\hat{g}_n - g\|_{w,2}^2$  where

Table 1.  $g(\theta) = 1$

Sample size $n$	Truncation parameter $m_n$				
	0	1	2	3	4
1000	0.251	0.128	0.330	5.186	110.752
10 000	0.250	0.113	0.089	0.555	11.100
100 000	0.250	0.111	0.065	0.091	1.135

Table 2.  $g(\theta) = (\pi/2) \sin(\pi\theta)$

Sample size $n$	Truncation parameter $m_n$				
	0	1	2	3	4
1000	0.48445	0.02041	0.32504	6.16375	128.20208
10 000	0.48378	0.00333	0.03352	0.61638	12.82021
100 000	0.48371	0.00162	0.00437	0.06164	1.28202

Table 3.  $g(\theta) = \exp(\theta)/(e - 1)$

Sample size $n$	Truncation parameter $m_n$				
	0	1	2	3	4
1000	0.558	0.291	0.431	5.701	126.007
10 000	0.558	0.277	0.184	0.660	12.663
100 000	0.558	0.275	0.159	0.156	1.329

Table 4.  $g(\theta) = \exp(-\theta)$ 

Sample size $n$	Truncation parameter $m_n$						
	0	1	2	3	4	5	6
500	0.08383	0.02271	0.01030	0.01426	0.03335	0.08570	0.22609
1000	0.08358	0.02177	0.00775	0.00778	0.01684	0.04289	0.11305
10 000	0.08336	0.02093	0.00546	0.00195	0.00198	0.00436	0.01132
100 000	0.08334	0.02084	0.00523	0.00137	0.00049	0.00051	0.00115

$w(\theta) = e^{-\theta}$ ,  $\forall \theta > 0$ . In this case  $\{p_{w_0, j}\}_{j=0}^{\infty}$  corresponds to the Laguerre polynomials with weight function  $w_0(\theta) = e^{-\theta}$ . Table 4 gives the integrated mean squared error of the estimator  $\hat{g}_{OP, n}$  when  $g(\theta) = e^{-\theta}$ ,  $\forall \theta > 0$  for sample sizes  $n = 1000, 10\ 000$  and  $100\ 000$  as well as for  $0 \leq m_n \leq 6$ .

*Remark 7.* Examples 1 and 2 (plus other unreported ones) indicate that  $\hat{g}_{OP, n}$  perform well for sample sizes  $n \geq 1000$  as long as  $h$ , defined as in (16), can be reasonably approximated by a low degree polynomial and that the optimal truncation parameter is used.

### 5. Estimating the optimal truncation parameter

In this section, a way is proposed to estimate the optimal truncation parameter  $m_n^*$  for the orthogonal polynomial mixing density estimator  $\hat{g}_{OP, n}$ , as given in (18), where  $m_n^*$  is defined to be the value of the truncation parameter  $m_n$  which minimizes  $E_g \|\hat{g}_{OP, n} - g\|_{w, 2}$ . We write

$$t_{n, j} = n^{-1} \{E_g \lambda_{w_0, j}^2(X_1) - (n+1)[E_g \lambda_{w_0, j}(X_1)]^2\}. \quad (30)$$

We observe from (29) that  $\sum_{j=0}^m t_{n, j} \leq \sum_{j=0}^m t_{n, j}$ , for all  $m \geq 0$ . This implies that  $m_n^*$  can be determined if the sign of  $\sum_{j=a}^b t_{n, j}$  is known for each  $a \leq b$ . Let  $\hat{t}_{n, j}$  be the unbiased estimator of  $t_{n, j}$  based on  $X_1, \dots, X_n$ ,  $\hat{t}_{n, i, j} = \sum_{l=i}^j \hat{t}_{n, l}$ ,  $\forall 0 \leq i \leq j$  and  $\sigma^2(\hat{t}_{n, i, j})$  be the unbiased estimator of the variance of  $\hat{t}_{n, i, j}$ . Let  $0 < \alpha^* < 1$  and  $B_n$  be the largest possible constant satisfying the inequalities

$$\max_{0 \leq x \leq j \leq B_n} \log(|k_{w_0, j, x}|/q(x)) \leq \beta_0 \log n, \quad B_n \leq \beta_1 \log n, \quad (31)$$

for positive constants  $\beta_0 < 1/2$  and  $\beta_1$ . Our algorithm for estimating  $m_n^*$  is as follows:

*Step 1.* Set  $\hat{m}_n^* = 0$  and  $n_1 = n_2 = 1$ .

*Step 2.* Compute  $Y = \hat{t}_{n, n_1, n_2} + z_{\alpha^*} \hat{\sigma}(\hat{t}_{n, n_1, n_2})$ , where  $\Phi(z_{\alpha^*}) = 1 - \alpha^*$  and  $\Phi$  denotes the distribution function of the standard normal distribution.

*Step 3.* If  $Y < 0$  and  $n_2 \leq B_n$ , then set  $\hat{m}_n^* = n_2$ ,  $n_1 = n_2 + 1$  and then set  $n_2 = n_1$ . Let  $Y = \hat{t}_{n, n_1, n_2} + z_{\alpha^*} \hat{\sigma}(\hat{t}_{n, n_1, n_2})$  and return to the start of Step 3.

*Elseif*  $Y \geq 0$  and  $n_2 \leq B_n$ , then increase  $n_2$  by 1, compute  $Y = \hat{t}_{n, n_1, n_2} + z_{\alpha^*} \hat{\sigma}(\hat{t}_{n, n_1, n_2})$  and return to the beginning of Step 3.

*Elseif*  $n_2 > B_n$ , then the estimate of the optimal truncation parameter  $m_n^*$  is given by  $\hat{m}_n^*$ .

*Endif.*

*Remark 8.* The above algorithm can be thought of as a successive sequence of hypothesis tests each at level  $\alpha^*$  where the null hypothesis always has fewer terms than the alternative.

*Remark 9.* The constant  $B_n$  can be chosen in the following manner. Under the conditions of corollary 1, take  $B_n = \beta_0 \log n / \log(2\gamma_1)$ . Under the conditions of corollary 2(a) and (b), we take  $B_n = \beta_0 \log n / \log(\gamma_1 \zeta)$  and  $(\beta_0/\gamma) \log n / \log \log n$  respectively.

*Remark 10.* The closer  $\alpha^*$  is chosen to 0, the more likely it is that  $\hat{m}_n^*$  will underestimate  $m_n^*$ . The previous section (see Tables 1 to 4) indicates that the risk of  $\hat{g}_{OP,n}$  is asymmetrical about  $m_n^*$  and that there is a distinct possibility that the risk increases very dramatically with overestimation. As such we recommend that  $\alpha^*$  be chosen to be 0.01, 0.05 or 0.10, which are in line with the usual values of  $\alpha^*$  for classical hypothesis testing.

Let  $\hat{g}_{OP,n}^*$  be as in (18) with  $m_n$  replaced by  $\hat{m}_n^*$ . The following theorem gives an upper bound to the convergence rate of  $\hat{g}_{OP,n}^*$ .

**Theorem 6**

Let  $\|C^2/w\|_1 < \infty$  and  $B_n$  be the largest possible constant satisfying (31). Then

$$\sup \{E_g \|\hat{g}_{OP,n}^* - g\|_{w,2} : g \in \mathcal{G}(\alpha, m, M, w_0)\} = O(B_n^{-\alpha}).$$

*Remark 11.* By choosing  $B_n = m_n$  in corollaries 1 and 2, we observe that the estimators  $\hat{g}_{OP,n}^*$  essentially retain the convergence rates of  $\hat{g}_{OP,n}$ .

*Example 1 continued.* Here we have applied the above algorithm to example 1. In particular the algorithm is used to determine  $\hat{m}_n^*$  using  $\alpha^* = 0.05$  and  $B_n = \lfloor \frac{1}{2} \log n \rfloor$ . For convenience we use 50 bootstrap replications to approximate each  $\hat{\sigma}(\hat{t}_{n,n_1,n_2})$ . The second

Table 5.  $g(\theta) = 1$

Sample size $n$	IMSE	Relative frequency of $\hat{m}_n^*$				
		0	1	2	3	4
1000	0.231	0.84	0.16	0.00	0.00	0.00
10 000	0.110	0.00	0.98	0.02	0.00	0.00
100 000	0.072	0.00	0.29	0.71	0.00	0.00

Table 6.  $g(\theta) = (\pi/2) \sin(\pi\theta)$

Sample size $n$	IMSE	Relative frequency of $\hat{m}_n^*$				
		0	1	2	3	4
1000	0.0552	0.08	0.92	0.00	0.00	0.00
10 000	0.00320	0.00	1.00	0.00	0.00	0.00
100 000	0.00158	0.00	1.00	0.00	0.00	0.00

Table 7.  $g(\theta) = \exp(\theta)/(e - 1)$

Sample size $n$	IMSE	Relative frequency of $\hat{m}_n^*$				
		0	1	2	3	4
1000	0.360	0.30	0.70	0.00	0.00	0.00
10 000	0.263	0.00	0.94	0.06	0.00	0.00
100 000	0.142	0.00	0.00	1.00	0.00	0.00

column of Tables 5, 6 and 7 give the average value of

$$(1/10) \sum_{i=1}^{10} [\hat{g}_{OP,n}^*(0.1i - 0.05) - g(0.1i - 0.05)]^2,$$

for 100 independent replications of  $X_1, \dots, X_n$ . These values approximate the integrated mean squared error (IMSE) of the mixing density estimator

$$\hat{g}_{OP,n}^*(\theta) = [C(\theta)/w(\theta)] \sum_{j=0}^{\hat{m}_n^*} \hat{h}_{w_0,j} P_{w_0,j}(\theta), \quad \forall 0 < \theta < \theta^*.$$

We recall that in this case, we have  $\theta^* = 1$ ,  $w(\theta) = 1$  and  $C(\theta) = 1 - \theta$ .

*Example 2 continued.* The above algorithm is now applied in example 2 with  $\alpha^* = 0.05$ ,  $B_n = 0.7 \log n$  and 50 bootstrap replications to approximate each  $\hat{\sigma}(\hat{t}_{n,n_1,n_2})$ . As in example 1, the second column of Table 8 gives the average value of

$$(1/10) \sum_{i=1}^{500} \exp \left\{ -(0.1i - 0.05) [\hat{g}_{OP,n}^*(0.1i - 0.05) - g(0.1i - 0.05)]^2 \right\},$$

for 100 independent replications of  $X_1, \dots, X_n$ . These values approximate the integrated weighted mean squared error (IMSE) of the orthogonal polynomial mixing density estimator  $\hat{g}_{OP,n}^*$ , namely  $E_g \|\hat{g}_{OP,n}^* - g\|_{w,2}^2$  with  $w(\theta) = e^{-\theta}$ ,  $\forall \theta > 0$ .

Table 8.  $g(\theta) = \exp(-\theta)$

Sample size $n$	IMSE	Relative frequency of $\hat{m}_n^*$				
		0	1	2	3	4
1000	0.0190	0.00	0.75	0.24	0.01	0.00
10 000	0.00455	0.00	0.00	0.71	0.29	0.00
100 000	0.00111	0.00	0.00	0.00	0.66	0.34

Both of the above Monte Carlo studies indicate that the risks of the orthogonal polynomial mixing density estimators  $\hat{g}_{OP,n}^*$  compare well to the ones based on the optimal truncation parameter.

We conclude with the remark that in general the following two conditions do not hold:  $\hat{g}_{OP,n}^*(\theta) \geq 0$ ,  $\forall 0 < \theta < \theta^*$  and  $\int_0^{\theta^*} \hat{g}_{OP,n}^*(\theta) d\theta = 1$ . As such the accuracy of estimate  $\hat{g}_{OP,n}^*$  can be further gauged by how close the above two conditions are to being satisfied.

### Acknowledgements

The research of Wei-Liem Loh was supported in part by NSA Grant MDA 904-93-3011 and the research of Cun-Hui Zhang was supported in part by ARO Grant DAAL03-91-G-0045 and NSA Grant R504-9282. We would like to thank Nicholas Hengartner for sending us a preprint of his paper.

### References

- Carroll, R. J. & Hall, P. (1988). Optimal rates of convergence for deconvolving a density. *J. Amer. Statist. Assoc.* **83**, 1184–1186.  
 Datta, S. (1991). On the consistency of posterior mixtures and its applications. *Ann. Statist.* **19**, 338–353.  
 DeVore, R. A. & Lorentz, G. G. (1993). *Constructive approximation*. Springer, New York.

- Devroye, L. P. & Wise, G. L. (1979). On the recovery of discrete probability densities from imperfect measurements. *J. Franklin Inst.* **307**, 1–20.
- Fan, J. (1991). On the optimal rates of convergence for nonparametric deconvolution problems. *Ann. Statist.* **19**, 1257–1272.
- Hengartner, N. W. (1995). Adaptive demixing in Poisson mixture models. Preprint.
- Johnstone, I. M. & Silverman, B. W. (1990). Speed of estimation in positron emission tomography and related inverse problems. *Ann. Statist.* **18**, 251–280.
- Loh, W. L. & Zhang, C.-H. (1996). Global properties of kernel estimators for mixing densities in exponential family models for discrete variables. *Statist. Sinica* **6**, 561–578.
- Meeden, G. (1972). Bayes estimation of the mixing distribution, the discrete case. *Ann. Math. Statist.* **43**, 1993–1999.
- Rolph, J. E. (1968). Bayesian estimation of mixing distributions. *Ann. Math. Statist.* **39**, 1289–1302.
- Simar, L. (1976). Maximum likelihood estimation of a compound Poisson process. *Ann. Statist.* **4**, 1200–1209.
- Szegő, G. (1975). *Orthogonal polynomials*. American Mathematical Society, Providence, RI.
- Tucker, H. G. (1963). An estimate of the compounding distribution of a compound Poisson distribution. *Theoret. Probab. Appl.* **8**, 195–200.
- Walter, G. G. & Hamedani, G. G. (1989). Bayes empirical Bayes estimation for discrete exponential families. *Ann. Inst. Statist. Math.* **41**, 101–119.
- Walter, G. G. & Hamedani, G. G. (1991). Bayes empirical Bayes estimation for natural exponential families with quadratic variance functions. *Ann. Statist.* **19**, 1191–1224.
- Zhang, C.-H. (1990). Fourier methods for estimating mixing densities and distributions. *Ann. Statist.* **18**, 806–831.
- Zhang, C.-H. (1995). On estimating mixing densities in discrete exponential family models. *Ann. Statist.* **23**, 929–945.

Received November 1994, in final form March 1996

Wei-Liem Loh, Department of Mathematics, National University of Singapore, Kent Ridge, Singapore 119260, Republic of Singapore.

Cun-Hui Zhang, Department of Statistics, Hill Center, Busch Campus, Rutgers University, New Brunswick, New Jersey 08903, USA.

## Appendix

*Proof of proposition 1.* We observe that

$$\begin{aligned} E_g \int_0^{\theta^*} [\hat{g}_{OP,n}^*(\theta) - g(\theta)]^2 w(\theta) d\theta &= E_g \int_0^{\theta^*} \left[ \sum_{j=0}^{m_n} \hat{h}_{w_0,j} p_{w_0,j}(\theta) - h(\theta) \right]^2 w_0(\theta) d\theta \\ &= \sum_{j=0}^{m_n} E_g (\hat{h}_{w_0,j} - h_{w_0,j})^2 + \sum_{j=m_n+1}^{\infty} h_{w_0,j}^2. \end{aligned} \quad (32)$$

The last equality follows from the completeness of  $\{p_{w_0,j}\}_{j=0}^{\infty}$ . Since

$$\hat{h}_{w_0,j} = n^{-1} \sum_{i=1}^n \lambda_{w_0,j}(X_i),$$

the right hand side of (32) is bounded by

$$\begin{aligned} n^{-1} \sum_{j=0}^{m_n} E_g [\lambda_{w_0,j}^2(X_1)] + \sum_{j=m_n+1}^{\infty} h_{w_0,j}^2 \\ \leq n^{-1} \sum_{j=0}^{m_n} \max_{0 \leq x \leq j} [k_{w_0,j,x}/q(x)]^2 + \sum_{j=m_n+1}^{\infty} h_{w_0,j}^2. \end{aligned}$$

This proves the proposition.

*Proof of proposition 2.* We observe from (19), (21) and repeated integration by parts that

$$\begin{aligned} \int_0^{\theta^*} h^{(m)}(\theta) p_{w_1, j}(\theta) w_1(\theta) d\theta &= (-1)^m \int_0^{\theta^*} h(\theta) [p_{w_1, j}(\theta) w_1(\theta)]^{(m)} d\theta \\ &= v_{j+m, m} \int_0^{\theta^*} h(\theta) p_{w_0, j+m}(\theta) w_0(\theta) d\theta \\ &= v_{j+m, m} h_{w_0, j+m} \quad \forall j \geq 0. \end{aligned}$$

From the completeness of  $\{p_{w_1, j}\}_{j=0}^{\infty}$ , we get  $\|h^{(m)}\|_{w_1, 2}^2 = \sum_{j=m}^{\infty} v_{j, m}^2 h_{w_0, j}^2$ . Now (22) follows immediately from (20).

*Proof of remark 3.* We begin with Laguerre polynomials.

*Laguerre polynomials.* Suppose  $w_0(\theta) = \theta^\beta e^{-\theta}$ , with  $\theta > 0$  and  $\beta > -1$ , is the weight function of the normalized Laguerre polynomials

$$p_{w_0, j}(\theta) = \left[ \Gamma(\beta + 1) \binom{j + \beta}{j} \right]^{-1/2} \sum_{x=0}^j \binom{j + \beta}{j - x} \frac{(-\theta)^x}{x!}, \quad \forall j \geq 0.$$

For  $j \geq 0$  and  $m \geq 1$ , we write  $w_1(\theta) = \theta^{\beta+m} e^{-\theta}$ ,

$$p_{w_1, j}(\theta) = \left[ \Gamma(\beta + m + 1) \binom{j + \beta + m}{j} \right]^{-1/2} \sum_{x=0}^j \binom{j + \beta + m}{j - x} \frac{(-\theta)^x}{x!},$$

and

$$v_{j+m, m} = \frac{(j+m)!}{j!} \left[ \Gamma(\beta + 1) \binom{j + \beta + m}{j+m} \right]^{1/2} \left[ \Gamma(\beta + m + 1) \binom{j + \beta + m}{j} \right]^{-1/2}.$$

Then (19) follows from the Rodrigues' formula for Laguerre polynomials and (20) holds for  $\alpha = m/2$ .

*Jacobi polynomials.* Suppose  $w_0(\theta) = \theta^{\beta_1} (\theta^* - \theta)^{\beta_2}$ , with  $\beta_1 > -1$ ,  $\beta_2 > -1$  and  $0 < \theta < \theta^* < \infty$ . Then the orthogonal polynomials with  $w_0$  as the weight function correspond to the normalized Jacobi polynomials

$$p_{w_0, j}(\theta) = C_{j, \beta_1, \beta_2} \binom{j + \beta_2}{j} (\theta^*)^{-j} \sum_{x=0}^j \frac{j(j-1) \cdots (j-x+1)}{(\beta_2 + 1)(\beta_2 + 2) \cdots (\beta_2 + x)} \binom{j + \beta_1}{x} \theta^{j-x} (\theta - \theta^*)^x,$$

where

$$C_{j, \beta_1, \beta_2} = \left[ \frac{(2j + \beta_1 + \beta_2 + 1) \Gamma(j+1) \Gamma(j + \beta_1 + \beta_2 + 1)}{(\theta^*)^{\beta_1 + \beta_2 + 1} \Gamma(j + \beta_1 + 1) \Gamma(j + \beta_2 + 1)} \right]^{1/2} \quad \text{if } j \geq 1,$$

and is equal to

$$\left[ \frac{\Gamma(\beta_1 + \beta_2 + 2)}{(\theta^*)^{\beta_1 + \beta_2 + 1} \Gamma(\beta_1 + 1) \Gamma(\beta_2 + 1)} \right]^{1/2} \quad \text{if } j = 0.$$

For  $m \geq 1$ , let  $p_{w_1, j}, j \geq 0$ , denote the set of normalized Jacobi polynomials with weight function

$$w_1(\theta) = \theta^{\beta_1 + m} (\theta^* - \theta)^{\beta_2 + m}, \quad \forall 0 < \theta < \theta^*,$$

and

$$v_{j+m, m} = (\theta^*)^m (j+m)! C_{j, \beta_1 + m, \beta_2 + m} / [j! C_{j+m, \beta_1, \beta_2}].$$

Then (19) follows from the Rodrigues' formula for Jacobi polynomials and (20) holds for  $\alpha = m$ .

*Proof of theorem 2.* We first observe from (23) that

$$n^{-1} \sum_{j=0}^{m_n} \max_{0 \leq x \leq j} [k_{w_0, j, x} / q(x)]^2 = O(m_n n^{2\beta_0 - 1}). \quad (33)$$

We also observe that

$$\sup \left\{ \sum_{j=m_n+1}^{\infty} h_{w_0, j}^2 : g \in \mathcal{G}(\alpha, m, M, w_0) \right\} = O(m_n^{-2\alpha}). \quad (34)$$

Now the theorem follows from (33), (34) and proposition 1.

*Proof of corollary 1.* From the properties of Laguerre polynomials, we have

$$\begin{aligned} |k_{w_0, j, x} / q(x)| &\leq \gamma_0 \gamma_1^x \binom{j+\beta}{j-x} \left[ \Gamma(\beta+1) \binom{j+\beta}{j} \right]^{-1/2} \\ &= \gamma_0 \gamma_1^x \binom{j}{x} \left[ \prod_{i=x+1}^j (1+\beta i^{-1}) \right]^{1/2} \left[ \Gamma(\beta+1) \prod_{i=1}^x (1+\beta i^{-1}) \right]^{-1/2} \\ &\leq \gamma_0 \gamma_1^j 2^j \left[ \prod_{i=x+1}^j (1+\beta i^{-1}) \right]^{1/2} \left[ \Gamma(\beta+1) \prod_{i=1}^x (1+\beta i^{-1}) \right]^{-1/2} \end{aligned} \quad (35)$$

Here we follow the convention that  $\prod_{i=x_1}^{x_2} (1+\beta i^{-1}) = 1$  if  $x_1 > x_2$ . We further observe that there exist positive constants  $c_1^*$  and  $c_2^*$  such that

$$c_1^* j^\beta \leq \prod_{i=1}^j (1+\beta i^{-1}) \leq c_2^* j^\beta, \quad \forall j \geq 1.$$

Thus it follows from (35) that

$$\max_{0 \leq x \leq j \leq m_n} \log (|k_{w_0, j, x} / q(x)|) = m_n (1 + o(1)) \log (2\gamma_1) \leq \beta_0 (1 + o(1)) \log n.$$

This proves (23) and the corollary follows from theorem 2.

*Proof of corollary 2.* If  $\gamma = 0$ , we observe that

$$\max_{0 \leq x \leq m_n} \log (1/q(x)) + m_n \log \zeta \leq m_n (1 + o(1)) \log (\gamma_1 \zeta) \leq \beta_0 (1 + o(1)) \log n.$$

This proves (25) and (a) follows from theorem 3. The case of  $0 < \gamma < \infty$  is similar and is omitted.

*Proof of theorem 4.* Let  $0 < \theta_0 < \theta_1 < a < \theta_2 < \theta_3 < \theta^*$  be fixed constants such that  $w$  is strictly positive and  $m$  times continuously differentiable on  $[\theta_0, \theta_3]$ . Define

$$h_{u, v}(\theta) = v^u \theta^{u-1} \exp(-v\theta) \Gamma(u),$$

$$g_{u, v}(\theta) = \{\chi_0(\theta) l_{1, u, v}(\theta) + \chi_1(\theta) h_{u, v}(\theta) + \chi_2(\theta) l_{2, u, v}(\theta)\} / C(\theta),$$

where  $\chi_j(\theta) = I\{\theta_j \leq \theta < \theta_{j+1}\}$  and  $l_{j, u, v}$ ,  $j = 1, 2$ , are polynomials each of degree  $(2m+1)$  such that  $g_{u, v}$  is  $m$  times continuously differentiable. Let  $g_0$  be an interior point of

$\mathcal{G}(\alpha, m, M, w_0)$ . Define

$$g_{0n}(\theta) = g_0(\theta) + \frac{3\varepsilon}{u_n^{1/4}} \left(\frac{\theta_2}{u_n}\right)^m \{g_{u_n, v_n}(\theta) + g_{00}(\theta) - (w_{0n} + 1)g_0(\theta)\},$$

$$g_{1n}(\theta) = g_{0n}(\theta) + \frac{\varepsilon}{u_n^{1/4}} \left(\frac{\theta_2}{u_n}\right)^m \left[ \sin\left(u_n \frac{\theta - a}{\theta_2}\right) - \frac{w_{1n}}{w_{0n}} \right] g_{u_n, v_n}(\theta),$$

$$g_{2n}(\theta) = g_{0n}(\theta) + \frac{\varepsilon}{u_n^{1/4}} \left(\frac{\theta_2}{u_n}\right)^m \left[ \cos\left(u_n \frac{\theta - a}{\theta_2}\right) - \frac{w_{2n}}{w_{0n}} \right] g_{u_n, v_n}(\theta),$$

where  $g_{00}$  is a density in  $\mathcal{G}(\alpha, m, M, w_0)$  bounded away from 0 in  $[\theta_0, \theta_3]$ ,  $w_{jn}$  are constants determined by  $\int_0^{\theta^*} g_{jn}(\theta) d\theta = 1$ ,  $\varepsilon > 0$ ,  $u_n = \delta_0 \log n$ , and  $v_n = u_n/a$ , with

$$\delta_0 = \max \left\{ \frac{\theta_2/(\theta_3 - \theta_2)}{\log(\theta_3/\theta_2)}, \frac{2}{\log(1 + a^2/\theta_2^2)}, \frac{1}{\theta_1/a - 1 - \log(\theta_1/a)}, \frac{1}{\theta_2/a - 1 - \log \theta_2/a} \right\}.$$

The rest of the proof is almost identical to steps 1 to 3 of the proof of th. 3 of Loh & Zhang (1996). It remains only to verify that  $g_{jn} \in \mathcal{G}(\alpha, m, M, w_0)$  for  $j = 0, 1, 2$ . Motivated by (16), define for  $0 < \theta < \theta^*$ ,

$$h(\theta) = \frac{3\varepsilon}{u_n^{1/4}} \left(\frac{\theta_2}{u_n}\right)^m w(\theta) g_{u, v}(\theta) / C(\theta).$$

Then using Leibniz rule we have  $\|h^{(m)}\|_{w_{1,2}} = \varepsilon O(1)$ , where the  $O(1)$  term does not depend on  $\varepsilon$ . Since (19) and (20) hold, we observe from proposition 2 that  $(\sum_{j=m}^{\infty} j^{2\alpha} h_{w_0, j}^2)^{1/2} = \varepsilon O(1)$ , where  $h_{w_0, j} = \int_0^{\theta^*} h(\theta) p_{w_0, j}(\theta) d\theta$ . Writing

$$g_{0n, w_0, j} = \int_0^{\theta^*} C(\theta) g_{0n}(\theta) p_{w_0, j}(\theta) d\theta, \quad \forall j \geq m,$$

it follows from Minkowski's inequality that  $(\sum_{j=m}^{\infty} j^{2\alpha} g_{0n, w_0, j}^2)^{1/2} < M + \varepsilon O(1)$  as  $g_0$  is an interior point of  $\mathcal{G}(\alpha, m, M, w_0)$ . Thus we conclude that  $g_{0n} \in \mathcal{G}(\alpha, m, M, w_0)$  for sufficiently small  $\varepsilon$ . Likewise we have  $g_{jn} \in \mathcal{G}(\alpha, m, M, w_0)$ ,  $j = 1, 2$ .

*Proof of theorem 5.* Let  $0 < \theta_1 < \theta_2 < \theta^*$  and  $\delta > 0$  be fixed constants such that  $w \geq \delta$  and  $m$  times continuously differentiable on  $[\theta_1, \theta_2]$  and  $g_0$  be an interior point of  $\mathcal{G}(\alpha, m, M, w_0)$  which is also bounded away from 0 on  $[\theta_1, \theta_2]$ . Furthermore we write

$$m_n = \lfloor \gamma \log n / \log \log n \rfloor, \quad (36)$$

where  $\gamma > 0$  is a suitably large constant to be determined later, and

$$T_n(\theta) = \sum_{j=0}^{m_n + 2m + 2} a_{n, j} \cos \left[ 2\pi j \left( \frac{\theta - \theta_1}{\theta_2 - \theta_1} \right) \right] I\{\theta_1 < \theta < \theta_2\}, \quad \forall 0 < \theta < \theta^*,$$

where the coefficients  $a_{n, j}$  are determined by the conditions that  $T_n$  is  $m$  times continuously differentiable on  $(0, \theta^*)$  and

$$\int_0^{\theta^*} T_n(\theta) d\theta = 0,$$

$$\int_0^{\theta^*} T_n(\theta) f(x | \theta) d\theta = 0, \quad \forall x = 0, 1, \dots, m_n - 1. \quad (37)$$



We observe from DeVore & Lorentz (1993) pp. 98–102 that

$$\|T_n^{(m)}\|_2 \leq \left[ \frac{2\pi(m_n + 2m + 2)}{\theta_2 - \theta_1} \right]^m \|T_n\|_2, \tag{38}$$

$$\|T_n\|_\infty \leq \left[ \frac{2m_n + 4m + 5}{\theta_2 - \theta_1} \right]^{1/2} \|T_n\|_2. \tag{39}$$

Define  $g_{n,1} = g_0$  and  $g_{n,2} = g_0 + \varepsilon_n T_n$  where  $\varepsilon_n = \varepsilon m_n^{-m} \|T_n\|_2^{-1}$  and  $\varepsilon > 0$  a suitably small constant (independent of  $n$ ) such that  $g_{n,2}$  is a density on  $(0, \theta^*)$ . Motivated by (16), define  $h(\theta) = \varepsilon_n w(\theta) T_n(\theta) / C(\theta)$  for  $0 < \theta < \theta^*$ . We observe from (38) that  $\|h^{(m)}\|_{w_1,2} = \varepsilon O(1)$ , where the  $O(1)$  term does not depend on  $\varepsilon$ . Hence arguing as in the proof of theorem 4, we conclude that  $g_{n,2} \in \mathcal{G}(\alpha, m, M, w_0)$  for sufficiently small  $\varepsilon$ .

Next we observe from (1), (27), (36) and (37) that for  $\lambda > 1$  and  $\gamma$  sufficiently large,

$$\begin{aligned} P_{g_{n,2}} \left[ \prod_{j=1}^n \frac{f(X_j; g_{n,2})}{f(X_j; g_{n,1})} \leq \lambda \right] &\geq 1 - P_{g_{n,2}} \left( \max_{1 \leq j \leq n} X_j \geq m_n \right) \\ &\geq 1 - O(1) n (\theta^* \tau_1)^{m_n} / (m_n)!^\tau \\ &= 1 + o(1). \end{aligned}$$

Finally we observe that  $\|g_{n,1} - g_{n,2}\|_{w,2} \geq (\varepsilon \delta^{1/2} / \gamma^m) (\log \log n / \log n)^m$  and now theorem 5 follows from lem. 1 of Zhang (1995).

*Proof of theorem 6.* Let  $g \in \mathcal{G}(\alpha, m, M, w_0)$  and  $h_{w_0,j}$  be as in (17). Define for each  $\beta > 0$ ,

$$j_n^*(\beta) = \begin{cases} \max \{j: 0 \leq j \leq B_n, h_{w_0,j}^2 > (\log n)^{-\beta}\} \\ \text{if } \{j: 0 \leq j \leq B_n, h_{w_0,j}^2 > (\log n)^{-\beta}\} \neq \emptyset, \\ 0 \end{cases} \quad \text{otherwise.}$$

We shall first show that

$$\sup \{P_g[\hat{m}_n^* < j_n^*(\beta)]: g \in \mathcal{G}(\alpha, m, M, w_0)\} = O(1)(\log n)^{2(1+\beta)n^{2\beta_0-1}}. \tag{40}$$

Since (40) clearly holds when  $j_n^*(\beta) = 0$ , it suffices to assume that  $j_n^*(\beta) \geq 1$ . Let  $t_{n,i,j} = \sum_{l=i}^j t_{n,l}$  and  $\sigma(\hat{t}_{n,i,j})$  be the standard deviation of  $\hat{t}_{n,i,j}$ . We observe from (30) and the definition of  $\lambda_{w_0,j}$  that  $\sup \{t_{n,j,j_n^*(\beta)}: g \in \mathcal{G}(\alpha, m, M, w_0), 0 \leq j \leq j_n^*(\beta)\} \leq -(\log n)^{-\beta/2}$  for sufficiently large  $n$ . Also

$$\begin{aligned} P_g[\hat{m}_n^* < j_n^*(\beta)] &\leq \sum_{j=0}^{j_n^*(\beta)-1} P_g[\hat{t}_{n,j+1,j_n^*(\beta)} + z_{\alpha^*}(\hat{t}_{n,j+1,j_n^*(\beta)}) \geq 0] \\ &= \sum_{j=0}^{j_n^*(\beta)-1} P_g \left[ \frac{\hat{t}_{n,j+1,j_n^*(\beta)} - t_{n,j+1,j_n^*(\beta)}}{\sigma(\hat{t}_{n,j+1,j_n^*(\beta)})} + z_{\alpha^*} \left( \frac{\hat{\sigma}(\hat{t}_{n,j+1,j_n^*(\beta)})}{\sigma(\hat{t}_{n,j+1,j_n^*(\beta)})} - 1 \right) \right] \\ &\geq -z_{\alpha^*} - \frac{t_{n,j+1,j_n^*(\beta)}}{\sigma(\hat{t}_{n,j+1,j_n^*(\beta)})} \\ &\leq 8(1+o(1))B_n(1+4z_{\alpha^*}^2)(\log n)^{2\beta} \sup \{\sigma^2(\hat{t}_{n,j+1,j_n^*(\beta)}): 0 \leq j < j_n^*(\beta)\}, \end{aligned} \tag{41}$$

uniformly over  $g \in \mathcal{G}(\alpha, m, M, w_0)$ . (40) now follows from (41) and the observation that

$$\sup \{\sigma^2(\hat{t}_{n,j+1,j_n^*(\beta)}): g \in \mathcal{G}(\alpha, m, M, w_0), 0 \leq j < j_n^*(\beta)\} = O((\log n)^{2n^{2\beta_0-1}}).$$

In a similar manner, we have

$$\sup \left\{ \sum_{j=1}^{m-1} h_{w_0, j}^2 P_g(\hat{m}_n^* < j) : g \in \mathcal{G}(\alpha, m, M, w_0) \right\} = o(B_n^{-2\alpha}). \quad (42)$$

Next as in (32), we observe that

$$\begin{aligned} & E_g \int_0^{\theta^*} [\hat{g}_{OP, n}^*(\theta) - g(\theta)]^2 w(\theta) d\theta \\ & \leq E_g \left\{ n^{-1} \sum_{j=0}^{B_n} \max_{0 \leq x \leq j} [k_{w_0, j, x}/q(x)]^2 + \sum_{j=B_n+1}^{\infty} h_{w_0, j}^2 \right. \\ & \quad \left. + \sum_{j=(\hat{m}_n^* + 1) \vee m}^{B_n} h_{w_0, j}^2 + \sum_{j=1}^{m-1} h_{w_0, j}^2 I\{\hat{m}_n^* < j\} \right\}. \end{aligned} \quad (43)$$

Conditioning on whether or not  $\hat{m}_n^* \geq j_n^*(\beta)$ , we observe using (40) that for sufficiently large  $\beta$ , the third term on the right hand side of (43) is bounded by

$$MP_g[\hat{m}_n^* < j_n^*(\beta)] + B_n(\log n)^{-\beta} = o(B_n^{-2\alpha}), \quad (44)$$

uniformly over  $g \in \mathcal{G}(\alpha, m, M, w_0)$  as  $n \rightarrow \infty$ . The theorem now follows from (42), (43) and (44).