# RATES OF CONVERGENCE FOR THE GAUSSIAN MIXTURE SIEVE

BY CHRISTOPHER R. GENOVESE[1] AND LARRY WASSERMAN[2]

*Carnegie Mellon University*

Gaussian mixtures provide a convenient method of density estimation that lies somewhere between parametric models and kernel density estimators. When the number of components of the mixture is allowed to increase as sample size increases, the model is called a mixture sieve. We establish a bound on the rate of convergence in Hellinger distance for density estimation using the Gaussian mixture sieve assuming that the true density is itself a mixture of Gaussians; the underlying mixing measure of the true density is not necessarily assumed to have finite support. Computing the rate involves some delicate calculations since the size of the sieve—as measured by bracketing entropy—and the saturation rate, cannot be found using standard methods. When the mixing measure has compact support, using $k_n \sim n^{2/3}/(\log n)^{1/3}$ components in the mixture yields a rate of order $(\log n)^{(1+\eta)/6}/n^{1/6}$ for every $\eta > 0$. The rates depend heavily on the tail behavior of the true density. The sensitivity to the tail behavior is diminished by using a robust sieve which includes a long-tailed component in the mixture. In the compact case, we obtain an improved rate of $(\log n/n)^{1/4}$. In the noncompact case, a spectrum of interesting rates arise depending on the thickness of the tails of the mixing measure.

**1. Introduction.** Statistical inference using mixtures of Gaussians is used for many purposes including density estimation, clustering and robust estimation; see, for example, Lindsay (1995), McLachlan and Basford (1988), Banfield and Raftery (1993) and Robert (1996). When the number of components of the mixture is allowed to increase with sample size, the model is called a Gaussian mixture sieve [Grenander (1981), Wong and Shen (1995)]. These sieves have been studied by several authors including Geman and Hwang (1982), Roeder (1992), Priebe (1994) and Roeder and Wasserman (1997). Related work from a Bayesian point of view is discussed in Escobar and West (1995). Priebe argues that in many cases, a mixture sieve has many advantages as a density estimator over kernel density estimates. For example, Priebe showed that with $n = 10,000$ observations, a log-normal density can be well approximated by a mixture of about 30 normals. In contrast, a kernel density estimator uses a mixture of 10,000 normals. Despite the ubiquity of mixture sieve models, little is known about their asymptotic properties. In particular, the rate of convergence of the density estimator of this sieve has not been established. In this paper, we bound the rate of convergence. Rates of

convergence for the mixing distribution function have been studied in Chen (1995). Also, van der Geer (1996) obtains rates for a different mixture model.

Let $\phi(x; \mu, \sigma)$ denote a Gaussian density with mean $\mu$ and variance $\sigma^2$. A finite Gaussian mixture is a density of the form

$$(1) \qquad f_\theta(x) = \sum_{j=1}^{k} p_j \phi(x; \mu_j, \sigma_j),$$

where $\theta = (\mu, \sigma, p)$, $\mu = (\mu_1, \ldots, \mu_k)$, $\sigma = (\sigma_1, \ldots, \sigma_k)$, $p = (p_1, \ldots, p_k)$. Here, the $\mu_j$'s are real, the $\sigma_j$'s are positive reals, $p_j \geq 0$ and $\sum_j p_j = 1$.

Let $m_k$, $s_k$ and $S$ be positive constants such that $m_k \to \infty$ and $s_k \downarrow 0$ as $k \to \infty$ and let

$$(2) \qquad \mathcal{M}_k = \left\{ f(\cdot) = \sum_{j=1}^{k} p_j \phi(\cdot\,; \mu_j, \sigma_j); |\mu_j| \leq m_k, \right.$$
$$\left. \text{and } s_k \leq \sigma_j \leq S, \; j = 1, \ldots, k \right\}.$$

Let $k_n$ be a sequence of integers such that $k_n \to \infty$ as $\infty$. The sieve we are interested in is $\mathcal{M}_{k_n}$. Our estimate of the true density is $\hat{f}_n(\cdot) = f_{\hat\theta(\cdot)}$ where $\hat\theta$ is the maximum likelihood estimate of $\theta$ in the model (2). We have chosen to fix $S$ mainly for convenience. This parameter can also be allowed to increase with $k$ but the results do not change materially.

We will assume that the true density is a "general Gaussian mixture" of the form

$$(3) \qquad f_0(x) = \int_0^\infty \int_{-\infty}^\infty \phi(x; \mu, \sigma) \, dP(\mu, \sigma)$$

for some probability measure $P$ on the Borel $\sigma$-algebra over $\mathcal{R} \times \mathcal{R}^+$. Let $\mathcal{F}$ denote all such densities. Of course, $\mathcal{F}$ contains all finite and countable mixtures as a special case. It is worth noting that the Dirichlet process mixture prior used in nonparametric Bayesian inference [Escobar and West (1995)] uses a prior with support in $\mathcal{F}$.

We measure the error in the estimate by Hellinger distance $d_H(f_0, \hat{f}_n)$ where $d_H(f, g)^2 = \int (\sqrt{f} - \sqrt{g})^2$. We bound the rate at which $d_H(f_0, \hat{f}_n)$ goes to 0, in two steps: (i) we bound the likelihood ratio outside Hellinger neighborhoods of the true density and (ii) we compute the rate at which finite mixtures saturate the set $\mathcal{F}$. The first task is addressed in Section 2 by computing the "size" of $\mathcal{M}_{k_n}$ in terms of its Hellinger bracketing entropy, and then appealing to recent results of Wong and Shen (1995). The second task is addressed in Section 3. Calculating the bracketing entropy and the saturation rate is usually straightforward for finite-dimensional models. However, mixture models do not behave as nicely as most finite-dimensional parametric families so these calculations require special attention. In particular, the square root of the density of a mixture model is not differentiable everywhere so that standard methods for computing entropy are not available. Hence, we believe that the

calculations in Sections 2 and 3 might be useful for other nonregular families as well. We put these pieces together and compute the rate in Section 4. Specifically we find $t_n$ such that $P^*(d_H(f_0, \hat{f}_n) > t_n) = o(1)$. In section 5 we discuss an improvement on the sieve. Section 6 contains closing remarks and unsolved problems.

The main conclusion of this paper is as follows. If the mixing measure $P$ is compactly supported, then taking $K_n \sim n^{2/3}(\log n)^{2/3}$ yields the rate $(\log n)^{(1+\eta)/6}/n^{1/6}$ for every $\eta > 0$. If the mixing measure is not compactly supported then we cannot compute the rate without the adjusted sieve in Section 5. In this case we get a rate of $(\log n)^{1/4}/n^{1/4}$ for the compact case and, in the noncompact case, we get a spectrum of rates depending on the tail behavior of $P$. We should mention that independently of our work, Li (1999) and Li and Barron (1999) also obtained a rate of $(\log n)^{1/4}/n^{1/4}$ for mixture models. More precisely, they obtained a rate in Kullback–Leibler distance, which corresponds to the above rate in Hellinger. Their proof is quite different from ours. On the one hand, it is more general since it applies to other mixtures besides Gaussian mixtures. On the other hand, their results do not directly apply to our case since their rates contain a constant which can be infinite and, furthermore, they assume the parameter space has been discretized. The results of Li and Barron are very interesting and they nicely complement the results in this paper.

REMARK. After a revision of this paper was submitted, Ghosal and van der Vaart (2000) obtained an improved rate of convergence for this problem. Our results are driven by the approximation error $\inf_{g \in \mathcal{M}_k} D(f_0, g) = O(\log k/k)$ where $D(f, g)$ is Kullback–Leibler distance. This implies that one needs $k(\varepsilon) \approx 1/\varepsilon$ mixture components to approximate an arbitrary $f_0$ to within $\varepsilon$ Kullback–Leibler distance. In bounding this approximation error we did not make use of the smoothness of the Gaussian densities. Ghosal and van der Vaart obtained an improved bound $\inf_{g \in \mathcal{M}_k} D(f_0, g) = O(\log k/e^k)$. This implies that one needs only $k(\varepsilon) \approx \log(1/\varepsilon)$ mixture components. As a consequence, they obtain a near parametric rate of $(\log n)^{\delta}/\sqrt{n}$ for some $\delta > 0$, in the case where the variances of the mixture components are bounded below by a known constant. This result appears to depend strongly on the smoothness of the Gaussians. Ghosal and van der Vaart did not obtain rates in the case where no such bound is known, though we believe that a $n^{-1/2}$ rate is not possible in this more general case. It appears that the $O(\log k/k)$ bound on the approximation error holds quite generally (i.e., without smoothness conditions on the densities being mixed) and could thus be used to obtain a convergence rate of $(\log n/n)^{1/4}$ without strong assumptions on the mixands. Indeed, Li and Barron (1999) obtained a bound of $O(1/k)$ with essentially no conditions on the mixands, though there are constants in their results which can be infinite in some cases. We suspect that these infinities can be eliminated at the expense of increasing the $O(1/k)$ term to $O(\log k/k)$. Currently, no results are available for the case where $m_k$, $s_k$ and $k$ are chosen using the data.

**2. Bounding the likelihood ratio.** In this section we bound the supremum of the likelihood outside a Hellinger neighborhood of the true density. Thoughout, we consider densities on the real line with respect to Lebesgue measure. Let $d_H(f, g)$ be Hellinger distance, $d_{TV}(f, g)$ total variation distance and $D(f, g)$ be Kullback–Leibler divergence; that is,

$$d_H^2(f, g) = \int \left(\sqrt{f}(x) - \sqrt{g}(x)\right)^2 dx = \int f(x)\,dx$$
$$+ \int g(x)\,dx - 2\int \sqrt{f(x)g(x)}\,dx,$$

$$d_{TV}(f, g) = \tfrac{1}{2}\int |f(x) - g(x)|\,dx,$$

$$D(f, g) = \int f(x)\log f(x)/g(x)\,dx.$$

The following inequalities are well known and will be used in what follows.

PROPOSITION. *Consider nonnegative, integrable functions f and g, not necessarily probability density functions and suppose that $g(x) \leq f(x)$ for almost all x with respect to Lebesgue measure. Then, $d_H(f, g) \leq \sqrt{2d_{TV}(f, g)}$ and $d_{TV}(f, g) \leq d_H(f, g)\sqrt{\int f(x)\,dx}$.*

2.1. *The bracketing entropy of $\mathscr{M}_k$.* In this subsection, we measure the size of $\mathscr{M}_k$ using bracketing entropy [van der Vaart and Wellner (1996), Section 2.7]. If $\mathscr{A}$ is a set of nonnegative, integrable functions and $d$ is a metric on this set, then an $\varepsilon$-bracketing (with respect to $d$) is a set of pairs of integrable functions $(l_1, u_1), \ldots, (l_m, u_m)$ such that (1) for each $f \in \mathscr{A}$ there exists $(l_j, u_j)$ such that $l_j \leq f \leq u_j$, a.e. with respect to Lebesgue measure and (2) $d(l_j, u_j) \leq \varepsilon$, $j = 1, \ldots, m$.

The smallest number of such brackets to cover $\mathscr{A}$ is called the *bracketing number* and is denoted by $N_{[]}(\varepsilon, \mathscr{A}, d)$. The *bracketing entropy* is defined by $H_{[]}(\varepsilon, \mathscr{A}, d) = \log N_{[]}(\varepsilon, \mathscr{A}, d)$.

Generally, if $\mathscr{A}$ is a parametric model of dimension $j$, then $N_{[]}(\varepsilon, \mathscr{A}, d) \sim \varepsilon^{-j}$ as can be proved using a Lipschitz argument; see, for example, van der Vaart and Wellner [(1996), Section 2.7.4]. But such arguments require that the derivative of the square root of the density be bounded by an $L_2$ function. This is not the case for mixtures. This is easy to see even in a simple mixture model like $(1-p)\phi(y; 0, 1) + p\phi(y; 0, 1/2)$; the derivative of the square root of this density at $p = 0$ behaves like $e^{x^2/2}$. Instead, we must bound the entropy by other methods. The result is given in the following theorem.

THEOREM 1. *Consider the set $\mathscr{M}_k$ defined by (2). If $\varepsilon \leq 1$, there exists positive constants $c_1$ and $c_2$, not depending on k or $\varepsilon$, such that*

$$N_{[]}(\varepsilon, \mathscr{M}_k, d_H) \leq c_1 c_2^k m_k^k \left(\frac{S}{s_k}\right)^{2k}\left(\frac{1}{\varepsilon}\right)^{3k-1}.$$

To prove Theorem 1, we need some lemmas.

LEMMA 1. *Let s, m and S be positive constants and define*

$$\mathscr{A} = \{\phi(\cdot; \mu, \sigma); |\mu| \le m, s \le \sigma \le S\}.$$

*Then, for $S \ge 1$ and $\varepsilon \in (0, 1)$,*

$$N_{[]}(\varepsilon, \mathscr{A}, d_H) \le \frac{128(2m)(S/s)^2}{\varepsilon^2}.$$

PROOF. Let $\delta = \varepsilon/2$ and $\tau^2 = (1 + \delta)S^2$. Let

$$r = \left\lceil \frac{4 \log \left(S\sqrt{1 + \delta}/s\right)}{\log \left(1 + \delta\right)} \right\rceil$$

where $\lceil a \rceil$ denotes the smallest integer greater than or equal to $a$. Define $\sigma_j^2 = \tau^2(1 + \delta)^{-j/2}$ for $j = 2, \ldots, r$. Note that $\sigma_r^2 \le s^2 \le S^2 = \sigma_2^2$. For $j \in \{2, \ldots, r\}$ let $\gamma_j = \delta\sigma_{j-2}/2$, let $I_j = \lceil m/\gamma_j \rceil$ and let $\mu_{ij} = i\gamma_j$ for $i = -I_j, \ldots, 0, \ldots, I_j$. Note that $[-m, m] \subset [-I_j\gamma_j, I_j\gamma_j]$. For $j \in \{2, \ldots, r\}$ and $-I_j \le i \le I_j$ let

(4) $$B_{ij} = \left\{(\mu, \sigma); \mu \in \left[\mu_{ij} - \frac{\delta\sigma_j}{4}, \mu_{ij} + \frac{\delta\sigma_j}{4}\right], \sigma^2 \in \left[\sigma_{j+1}^2, \sigma_j^2\right]\right\}.$$

The $B_{ij}$ cover the parameter space. [A similar construction is used in Tong and Viele (1998).]

Let

$$l_{ij}(y) = (1 + \delta)^{-1}\phi\left(y; \mu_{ij}, \frac{\sigma_{j+1}^2}{(1 + \delta)^{1/4}}\right)$$

and

$$u_{ij}(y) = (1 + \delta)\phi\left(y; \mu_{ij}, \sigma_j^2(1 + \delta)\right).$$

We claim that $(l_{ij}, u_{ij})$ brackets $B_{ij}$. This follows, after some algebra, from the fact that, whenever $\sigma_1 < \sigma_2$,

$$\frac{\phi(y; \mu_1, \sigma_1)}{\phi(y; \mu_2, \sigma_2)} \le \frac{\sigma_2}{\sigma_1} \exp\left\{\frac{(\mu_1 - \mu_2)^2}{2(\sigma_2^2 - \sigma_1^2)}\right\}.$$

Next we bound $d_H(l_{ij}, u_{ij})$. In general, if $f$ and $g$ are probability density functions then $d_H^2((1+\delta)f, (1+\delta)^{-1}g) = d_H^2(f, g) + \delta^2/(1+\delta) \le d_H^2(f, g) + \delta^2$. Also, if $f(y) = \phi(y; \mu_1, \sigma_1)$ and $g(y) = \phi(y; \mu_2, \sigma_2)$, then

$$d_H^2(f, g) = 2\left[1 - \left\{\frac{2\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2}\right\}^{1/2} \exp\left\{-\frac{(\mu_1 - \mu_2)^2}{4(\sigma_1^2 + \sigma_2^2)}\right\}\right].$$

So,

$$d_H^2(l_{ij}, u_{ij}) \le 2\left[1 - \sqrt{2}\left\{\frac{(1+\delta)^{7/8}}{(1+\delta)^{7/4}+1}\right\}^{1/2}\right] + \delta^2 \le 2\delta^2 \le \varepsilon^2.$$

The last line holds because of the following inequality:

$$1 - \left\{\frac{2u}{1+u^2}\right\}^{1/2} = \frac{(u-1)^2}{\left[(1+u^2)\left(1+\{2u/(1+u^2)\}^{1/2}\right)\right]} \le \frac{1}{2}(u-1)^2,$$

where $u = (1+\delta)^{7/8} > 1$ and $u - 1 \le \delta$.

Finally we count the number of boxes $N$. For each $j$, the number of boxes is less than or equal to $2m/\gamma_j$. Thus, we see that

$$N \le 2m \sum_{j=2}^{r} \frac{1}{\gamma_j} = \frac{4m}{\delta S\sqrt{1+\delta}} \sum_{j=2}^{r}(1+\delta)^{(j-2)/4} = \frac{4m}{\delta S\sqrt{1+\delta}} \sum_{j=2}^{r}(1+\delta)^{j/4}$$

$$\le \frac{4mr}{\delta S(1+\delta)}(1+\delta)^{r/4} = \frac{4mr}{\delta S(1+\delta)} \frac{S\sqrt{1+\delta}}{s} \le \frac{16m}{\delta s} \frac{\log\left(S\sqrt{1+\delta}/s\right)}{\log(1+\delta)}$$

$$\le \frac{32m}{\delta s} \frac{S\sqrt{1+\delta}}{\delta s} \le \frac{256mS}{\varepsilon^2 s^2} \le \frac{256mS^2}{\varepsilon^2 s^2}. \qquad \square$$

Let $\mathscr{S}_{k-1} = \{p = (p_1, \ldots, p_k); \ p_j \ge 0, \ \sum_j p_j = 1\}$ be the $k-1$ dimensional simplex.

LEMMA 2 (Bracketing entropy of the simplex). *If $\varepsilon \le 1$, then*

$$N_{[]}(\varepsilon, S_{k-1}, d_H) \le \frac{k(2\pi e)^{k/2}}{\varepsilon^{k-1}}.$$

PROOF. Given $p = (p_1, \ldots, p_k) \in \mathscr{S}_{k-1}$, let $q = (q_1, \ldots, q_k)$ where $q_j = \sqrt{p_j}$. Then $p \in \mathscr{S}_{k-1}$ if and only if $q \in Q^+ \cap U$ where $U$ is the surface of the unit sphere and $Q^+$ is the positive quadrant of $\mathscr{R}^k$. By virtue of this mapping, an $\varepsilon$-$L_2$ bracketing of $Q^+ \cap U$ corresponds to an $\varepsilon$-Hellinger bracketing of $\mathscr{S}_{k-1}$.

Divide the unit cube in $\mathscr{R}^k$ into disjoint cubes with sides parallel to the axes and with sides of length $\varepsilon/\sqrt{k}$. Let $\mathscr{C} = \{C_1, \ldots, C_N\}$ be the subset of these cubes that have non-empty intersection with $Q^+ \cap U$. Let $b_r$ be the vertex of $C_r$ furthest from the origin and let $a_r$ be the vertex of $C_r$ closest to the origin. Note that $\sum_j(a_{rj} - b_{rj})^2 = \varepsilon^2$ so $\{(a_1, b_1), \ldots, (a_N, b_N)\}$ forms an $\varepsilon$-$L_2$ bracketing. It remains to count the number of cubes $N$.

Let $T_a = \{q \in Q^+; \|q\| \le a\}$. Let $C = \bigcup_{C_j \in \mathscr{C}} C_j$. Note that $C \subset T_{1+\varepsilon} - T_{1-\varepsilon} \equiv A$ so

$$\text{Volume}(A) \ge \text{Volume}(C) = N\left(\frac{\varepsilon}{\sqrt{k}}\right)^k.$$

Let $V_k(a) = a^k \pi^{k/2}/\Gamma((k/2)+1)$ denote the volume of a sphere of radius $a$. Then,

$$N \leq \frac{\text{Volume }(A)}{\left(\varepsilon/\sqrt{k}\right)^k} = \frac{1}{2^k} \frac{[V_k(1+\varepsilon) - V_k(1-\varepsilon)]}{\left(\varepsilon/\sqrt{k}\right)^k}$$

$$= \frac{1}{2^k} \frac{[(1+\varepsilon)^k - (1-\varepsilon)^k]}{\left(\varepsilon/\sqrt{k}\right)^k} \frac{\pi^{k/2}}{(k/2)!} \leq \left(\frac{\pi e}{2}\right)^{k/2} \frac{[(1+\varepsilon)^k - (1-\varepsilon)^k]}{\varepsilon^k}$$

since $x! \geq x^x e^{-x}$. Now, $(1+\varepsilon)^k - (1-\varepsilon)^k = k \int_{(1-\varepsilon)}^{(1+\varepsilon)} x^{k-1}\, dx \leq 2\varepsilon k (1+\varepsilon)^{k-1}$. The conclusion follows. $\quad\square$

LEMMA 3.   *Let $(l_1, u_1), \ldots, (l_m, u_m)$ be any $\varepsilon$ Hellinger bracketing. If $\varepsilon \leq 1$ then*

$$\Delta_\varepsilon \equiv \max_j \int u_j(x)\, dx \leq 1 + 3\varepsilon.$$

PROOF.   Let $u$ denote one of the upper brackets and let $l$ denote the corresponding lower bracket. Then, $\int u = \|\sqrt{u}\|_2^2 \leq (\|\sqrt{l}\|_2 + \|\sqrt{u} - \sqrt{l}\|_2)^2 \leq (1+\varepsilon)^2 \leq 1 + 3\varepsilon$. $\quad\square$

LEMMA 4.   *Let $\{l_1, \ldots, l_k\}$ and $\{u_1, \ldots, u_k\}$ be nonnegative, integrable functions and let $(a_1, \ldots, a_k)$ and $(b_1, \ldots, b_k)$ be vectors of nonnegative real numbers. Let $l = \sum_{j=1}^k a_j l_j$ and $u = \sum_{j=1}^k b_j u_j$. Then,*

$$d_H^2(l, u) \leq \sum_{j=1}^k d_H^2(a_j l_j, b_j u_j).$$

PROOF.   Note that

$$d_H^2(l, u) = \sum_j b_j \int u_j + \sum_j a_j \int l_j - 2 \int \sqrt{\sum_j b_j u_j \sum_j a_j l_j}$$

and

$$\sum_{j=1}^k d_H^2(a_j l_j, b_j u_j) = \sum_j b_j \int u_j + \sum_j a_j \int l_j - 2 \int \sum_j \sqrt{a_j b_j l_j u_j}.$$

Thus, it suffices to show that

$$\int \sqrt{\sum_j b_j u_j \sum_j a_j l_j} \geq \int \sum_j \sqrt{a_j b_j l_j u_j}$$

and for this, it suffices to show that

$$(5) \qquad \sqrt{\sum_j b_j u_j \sum_j a_j l_j} \geq \sum_j \sqrt{a_j b_j l_j u_j}$$

for all $x$. This follows from the Cauchy–Schwartz inequality. $\square$

THEOREM 2. *Let $\mathscr{F}_j = \{f_{\theta_j}; \theta_j \in \Omega_j\}$ be a set of density functions for $j = 1, \ldots, k$. Let*

$$\mathscr{M}_k = \left\{ f(\cdot) = \sum_{j=1}^k p_j f_{\theta_j}(\cdot); \theta_j \in \Omega_j, p_j \geq 0, \sum_j p_j = 1 \right\}.$$

*If $\varepsilon \geq 1$ then*

$$(6) \qquad N_{[]}(\varepsilon, \mathscr{M}_k, d_H) \leq k(2\pi e)^{k/2} \left(\frac{3}{\varepsilon}\right)^{k-1} \prod_{j=1}^k N_{[]}(\varepsilon/3, \mathscr{F}_j, d_H).$$

PROOF. Let $\delta = \varepsilon/3$ and let

$$\mathscr{B}_j = \{(l_{j1}, u_{j1}), \ldots, (l_{jm}, u_{jm})\}$$

be a set of $\delta$ Hellinger brackets for $\mathscr{F}_j$. Let $\{(a_1, b_1), \ldots, (a_s, b_s)\}$ be a $\delta$ bracketing for the simplex $\mathscr{S}_{k-1}$. Note that each $a_r$ and $b_r$ is a vector of length $k$. From Lemma 3, $\max_j \int u_j(x)\,dx$ and $\max_r \sum_{j=1}^k b_{rj}$ are bounded above by $1 + 3\delta$, where $b_{rj}$ is the $j$th component of the vector $b_r$.

Consider $h(x) = \sum_j p_j f_{\theta_j}(x) \in \mathscr{M}_k$. Let $a = (a_1, \ldots, a_k)$ and $b = (b_1, \ldots, b_k)$ be an $\varepsilon$-bracket for $p = (p_1, \ldots, p_k)$. Let $(l_j, u_j) \in \mathscr{B}_j$ be an $\varepsilon$-bracket for $f_{\theta_j}$. Define $l = \sum_j a_j l_j$ and $u = \sum_j b_j u_j$. Clearly, $l(x) \leq h(x) \leq u(x)$ and the number of such brackets is bounded by the right-hand side of (6). Now we show that $d_H(l, u) \leq \varepsilon$.

By Lemma 4, $d_H^2(l, u) \leq \sum_{j=1}^k d_j^2$ where $d_j = d_H(b_j u_j, a_j l_j)$. Now, using the Cauchy–Schwarz inequality and the fact that $\int l_j \leq 1$ we have

$$d_j^2 = \int \left(\sqrt{a_j l_j} - \sqrt{b_j u_j}\right)^2 dx$$

$$= \int \left(\sqrt{b_j}\left(\sqrt{l_j} - \sqrt{u_j}\right) + \left(\sqrt{a_j} - \sqrt{b_j}\right)\sqrt{l_j}\right)^2 dx$$

$$= b_j \int \left(\sqrt{l_j} - \sqrt{u_j}\right)^2 dx + \left(\sqrt{a_j} - \sqrt{b_j}\right)^2 \int l_j\,dx$$

$$\quad + 2\sqrt{b_j}\left(\sqrt{b_j} - \sqrt{a_j}\right) \int \sqrt{l_j}\left(\sqrt{u_j} - \sqrt{l_j}\right)$$

$$\leq b_j \delta^2 + \left(\sqrt{a_j} - \sqrt{b_j}\right)^2 + 2\sqrt{b_j}\left(\sqrt{b_j} - \sqrt{a_j}\right)\delta.$$

Thus,

$$d_H^2(l, u) \leq \sum_j d_j^2 \leq \sum_j b_j \delta^2 + \sum_j \left(\sqrt{a_j} - \sqrt{b_j}\right)^2 + 2\delta \sum_j \sqrt{b_j}\left(\sqrt{b_j} - \sqrt{a_j}\right)$$

$$\leq \delta^2(1 + 3\delta) + \delta^2 + 2\delta\sqrt{(1 + 3\delta)}\delta \leq \delta^2[4 + 2\sqrt{3}] \leq \varepsilon^2. \qquad \square$$

The proof of Theorem 1 follows from Lemma 1, Lemma 2 and Theorem 2.

2.2. *Large deviation bound.* Now we use the results of the previous section to bound the likelihood ratio. We will need the following result from Wong and Shen (1995).

LEMMA 5 [Wong and Shen (1995), Theorem 1]. *Let* $X_1, \ldots, X_n$ *be i.i.d. from a distribution* $P_0$ *with density* $f_0$ *and define the likelihood ratio*

$$R_n(f) = \prod_{i=1}^n \frac{f(X_i)}{f_0(X_i)}.$$

*Let* $\mathscr{A}$ *be a set of density functions. There are positive constants* $c_1, c_2, c_3, c_4$ *such that if*

$$\int_{\varepsilon^2/2^8}^{\sqrt{2}\varepsilon} H_{[]}^{1/2}(u/c_3, \mathscr{A}, d_H)\, du \leq c_4 \sqrt{n}\varepsilon^2$$

*then*

$$P_0^*\left(\sup_{f \in B} R_n(f) > e^{-nc_1\varepsilon^2}\right) \leq 4e^{-c_2 n\varepsilon^2},$$

*where*

$$B = \{f \in \mathscr{A}; d_H(f_0, f) > \varepsilon\}.$$

THEOREM 3. *Let* $\mathscr{M}_{k_n}$ *be a mixture of* $k_n$ *Gaussians. Let* $X_1, \ldots, X_n$ *be i.i.d. from a distribution with some density* $f_0$. *Let* $\alpha, \alpha_0, \beta$ *and* $\beta_0$ *be nonnegative constants such that* $2\alpha + \beta \leq 1$ *and* $2\alpha_0 + \beta_0 \geq 1$. *Let* $k_n = n^\beta/(\log n)^{\beta_0}$ *and* $\varepsilon_n = (\log n)^{\alpha_0}/n^\alpha$. *Suppose that* $m_k/s_k = O(k^\eta)$ *for some* $\eta > 0$. *Then, with probability* 1, *there exists* $n_0, c_1$ *and* $c_2$ *such that, for all* $n \geq n_0$,

$$\sup_{f \in B_n} R_n(f) < c_1 e^{-c_2 n\varepsilon_n^2},$$

*where* $B_n = \{f \in \mathscr{M}_{k_n}; d_H(f_0, f) > \varepsilon_n\}$.

PROOF. It follows from Theorem 1 that $N_{[]}(\varepsilon_n, \mathscr{M}_{k_n}, d_H) \preceq (r/\varepsilon_n)^{3k_n}$ where $r \asymp (S^2 m_{k_n}/s_{k_n}^2)^{1/3}$. Let $a = \sqrt{2\log(r/(\varepsilon\sqrt{2}))}$ and $b = \sqrt{2\log(2^8 r/\varepsilon^2)}$ and use the substitution $u^2 = 2\log(r/x)$ to see that

$$
\begin{aligned}
J(\varepsilon) &\equiv \int_{\varepsilon_n^2/2^8}^{\sqrt{2}\varepsilon_n} H_{[]}^{1/2}(u/c_3, \mathscr{M}_{k_n}, d_H)\, du \preceq \sqrt{k_n} \int_{\varepsilon_n^2/2^8}^{\sqrt{2}\varepsilon_n} \sqrt{\log \frac{r}{x}}\, dx \\
&= r\sqrt{k_n/2} \int_a^b u^2 e^{-u^2/2}\, du \le r\sqrt{k_n/2} \int_a^\infty u^2 e^{-u^2/2}\, du \\
&= r\sqrt{k_n/2}\Big[ ae^{-a^2/2} + \int_a^\infty e^{-u^2/2}\, du \Big] \le r\sqrt{k_n/2}\Big[ ae^{-a^2/2} + \frac{\sqrt{2\pi}}{a}e^{-a^2/2} \Big] \\
&\le r\sqrt{k_n/2}\Big[ ae^{-a^2/2} + ae^{-a^2/2} \Big] \le 2\sqrt{2}\varepsilon_n \Big\{ k_n \log \frac{r}{\sqrt{2}\varepsilon_n} \Big\}^{1/2}.
\end{aligned}
$$

With $k_n$, $\varepsilon_n$ and $m_k/s_k$ chosen as in the statement of the theorem, it follows that $J(\varepsilon_n) \preceq \sqrt{n}\varepsilon_n^2$. The result follows from Lemma 5 and the first Borel–Cantelli lemma. □

**3. Saturation rate.** In this section we establish the saturation rate of the sieve in both Kullback–Leibler distance and $\chi^2$ distance. It turns out that the saturation rate depends on tail conditions on $f_0$, and mixtures turn out to require some special treatment. The usual saturation rate arguments used in function estimation theory do not readily apply.

Recall that the true density is assumed to be of the form

$$
f_0(y) = \int_0^\infty \int_{-\infty}^\infty \phi(y; \mu, \sigma)\, dP(\mu, \sigma).
$$

Let $m_k$, $s_k$ be sequences of positive real numbers with $s_k > 0$. Define

$$
R_k = \{(\mu, \sigma); |\mu| < m_k, s_k < \sigma < S\}
$$

and let $\delta_k = P(R_k^c)$. The following lemma will be useful.

LEMMA 6 [Barron and Yang (1995)]. *If $f/g \le V$ then*

$$
D(f, g) \le [2 + \log V]d_H^2(f, g).
$$

In what follows, it will be convenient to make a slight change in the definition of the sieve. Specifically, we now define $\mathscr{M}_k$ to be mixtures of $k+1$ components instead of $k$ components. This makes some bookkeeping in Theorem 4 simpler. Note that Theorem 3 is still true with this change.

THEOREM 4. *Let $m_k \to \infty$ and $s_k \to 0$. Let $r_k = \sqrt{8m_k/(ks_k)}$ and assume that $r_k = o(1)$ as $k \to \infty$ and that*

$$
E_P(\sigma^{-1}) < \infty.
$$

*Let $a_k$ be the smallest real number such that*

(7)
$$E_P\left(\frac{1}{\sigma}e^{\mu^2/a_k^2}\,\Big|\,R_k^c\right) \le a_k,$$

*where the expectation is taken to be the essential supremum over $R_k$ if $P(R_k^c) = 0$. Then, for $r_k < 2.5$,*

(8)
$$\inf_{f\in\mathcal{M}_k} D(f_0, f) \le 2\delta_k\Big[2 + \log\Big(1 + a_k\sqrt{S^2 + a_k^2}\Big)\Big]$$
$$+ \log(1 + \delta_k) + 2r_k \equiv \omega_k$$

PROOF.   Let
$$f_k(y) = \int_{s_k}^{S}\int_{-m_k}^{m_k} \phi(y;\mu,\sigma)\,dP(\mu,\sigma) + \delta_k\phi_0(y),$$

where $\phi_0(y)$ is a normal density with mean 0 and variance $S^2 + a_k^2$. Given a set of points $\{(\mu_1,\sigma_1),\ldots,(\mu_k,\sigma_k)\}$ and a partition $\{A_1,\ldots,A_k\}$ of $R_k$, to be chosen later, define

$$g_k = \sum_{j=1}^{k} p_j\phi(y;\mu_j,\sigma_j) + \delta_k\phi_0,$$

where $p_j = \int_{A_j} dP$. Let $h_k = g_k/(1 + \delta_k)$. Then $h_k \in \mathcal{M}_k$, $D(f_0, h_k) = D(f_0, g_k) + \log(1 + \delta_k)$ and

(9)
$$D(f_0, g_k) = D(f_0, f_k) + \int f_0 \log\frac{f_k}{g_k}.$$

To bound $D(f_0, f_k)$, it is helpful to first bound $f_0/f_k$. To this end, let

$$\gamma = \frac{\mu/\sigma^2}{1/\sigma^2 - 1/(S^2 + a_k^2)}$$

and note that

$$\frac{f_0}{f_k} = \frac{\int_{R_k}\phi\,dP}{\int_{R_k}\phi\,dP + \delta_k\phi_0} + \frac{\int_{R_k^c}\phi\,dP}{\int_{R_k}\phi\,dP + \delta_k\phi_0} \le 1 + \frac{1}{\delta_k}\int_{R_k^c}\frac{\phi}{\phi_0}\,dP$$

$$\le 1 + \frac{\sqrt{S^2 + a_k^2}}{\delta_k}\int_{R_k^c}\frac{1}{\sigma}\exp\left\{-\frac{1}{2}\Big[\frac{(x-\mu)^2}{\sigma^2} - \frac{x^2}{S^2 + a_k^2}\Big]\right\}dP$$

$$\le 1 + \frac{\sqrt{S^2 + a_k^2}}{\delta_k}\int_{R_k^c}\frac{1}{\sigma}\exp\left\{-\frac{1}{2}\Big[\frac{1}{\sigma^2} - \frac{1}{S^2 + a_k^2}\Big](x-\gamma)^2\right\}$$

$$\times \exp\left\{\frac{\mu^2}{\sigma^2(S^2 + a_k^2)(1/\sigma^2 - 1/(S^2 + a_k^2))}\right\}dP$$

$$\leq 1 + \frac{\sqrt{S^2 + a_k^2}}{\delta_k} \int_{R_k^c} \frac{1}{\sigma} \exp\left\{ \frac{\mu^2}{S^2 + a_k^2 - \sigma^2} \right\} dP$$

$$\leq 1 + \frac{\sqrt{S^2 + a_k^2}}{\delta_k} \int_{R_k^c} \frac{1}{\sigma} e^{\mu^2/a_k^2} dP \leq 1 + \sqrt{S^2 + a_k^2}\, a_k,$$

where the last inequality follows from (7). By Lemma 6, and the fact that $d_H^2(f_0, f_k) \leq 2 d_{TV}(f_0, f_k) \leq 2\delta_k$, we have

$$D(f_0, f_k) \leq 2\delta_k \left[ 2 + \log\left( 1 + a_k \sqrt{S^2 + a_k^2} \right) \right].$$

Thus we have bounded the first term in (9). Next we bound the second term.

Consider any $A_j$ in the partition. Let $\phi_j(x) = \phi(x; \mu_j, \sigma_j)$ and define

$$v^{-2} = \frac{1}{\sigma^2} - \frac{1}{\sigma_j^2} \quad \text{and} \quad \gamma = \frac{\mu/\sigma^2 - \mu_j/\sigma_j^2}{1/\sigma^2 - 1/\sigma_j^2}.$$

Then

$$\int_{A_j} \phi\, dP = \phi_j \int_{A_j} \frac{\phi}{\phi_j}\, dP = \phi_j \int_{A_j} \frac{\sigma_j}{\sigma} \exp\left( -\frac{1}{2} \left[ \frac{(x-\mu)^2}{\sigma^2} - \frac{(x-\mu_j)^2}{\sigma_j^2} \right] \right) dP$$

$$= \phi_j \int_{A_j} \frac{\sigma_j}{\sigma} \exp\left( -\frac{1}{2v^2}(x-\gamma)^2 \right) \exp\left( \frac{1}{2} \frac{(\mu-\mu_j)^2}{\sigma_j^2 - \sigma^2} \right) dP$$

$$\leq \phi_j \int_{A_j} \frac{\sigma_j}{\sigma} \exp\left( \frac{1}{2} \frac{(\mu-\mu_j)^2}{\sigma_j^2 - \sigma^2} \right) dP.$$

Below we will show that

(10) $$\left( \frac{\sigma_j}{\sigma} \right) \exp\left( \frac{1}{2} \frac{(\mu-\mu_j)^2}{\sigma_j^2 - \sigma^2} \right) \leq (1 + r_k)^2.$$

It then follows that

$$\int_{A_j} \phi\, dP \leq (1 + r_k)^2 p_j \phi_j.$$

Hence,

$$\frac{f_k}{g_k} = \frac{\sum_j \int_{A_j} \phi\, dP + r_k \phi_0}{\sum_j p_j \phi_j + r_k \phi_0} \leq \frac{(1 + r_k)^2 \sum_j p_j \phi_j + r_k \phi_0}{\sum_j p_j \phi_j + r_k \phi_0} \leq (1 + r_k)^2$$

and so

$$\int f_0 \log \frac{f_k}{g_k} \leq 2\log(r_k + 1) \leq 2r_k.$$

It remains to be shown that (10) holds and that $g_k \in \mathcal{M}_k$.

Let

$$v_1 < v_2 < \cdots < v_J \equiv S,$$

where

$$v_j = v_{j-1}\sqrt{1+r_k} \quad \text{and} \quad \sigma_j = v_j\sqrt{1+r_k}.$$

Here, $J$ is the largest integer such that $s_k \leq v_1$. Also, for $\sigma \in [v_{j-1}, v_j]$, divide $[-m_k, m_k]$ into intervals of length $\xi_j$ where

(11) $$\xi_j = v_j\sqrt{2r_k \log(1+r_k)}.$$

Then, for $\sigma \in [v_{j-1}, v_j]$, $\sigma_j/\sigma \leq \sigma_j/v_{j-1} = (1+r_k)$. Also, for $\sigma_j^2 - \sigma^2 \geq r_k v_j^2$. Hence,

$$\exp\left\{\frac{1}{2}\frac{(\mu - \mu_j)^2}{\sigma_j^2 - \sigma^2}\right\} \leq \exp\left\{\frac{1}{2}\frac{\xi_j^2}{r_k v_j^2}\right\} \leq (1+r_k)$$

using (11). Thus,

$$\frac{\sigma_j}{\sigma}\exp\left\{\frac{1}{2}\frac{(\mu - \mu_j)^2}{\sigma_j^2 - \sigma^2}\right\} \leq (1+r_k)^2$$

which confirms (10).

To ensure that $g_k \in \mathcal{M}_k$ we have to make sure that the above scheme partitions $R_k$ into no more than $k$ pieces. For fixed $\sigma_j$, the number of divisions of $\mu$ is

$$\frac{\text{total length}}{\xi_j} = \frac{2m_k}{v_j\sqrt{2r_k \log(1+r_k)}}.$$

So, using the fact that $r_k < 2.5$, the total number $N$ of rectangles is

$$N = \frac{2m_k}{\sqrt{2r_k \log(1+r_k)}}\left[\frac{1}{v_1} + \cdots + \frac{1}{v_J}\right]$$

$$= \frac{2m_k}{\sqrt{2r_k \log(1+r_k)}}\left[\frac{1}{s_k} + \frac{1}{s_k(1+r_k)^{1/2}} + \frac{1}{s_k(1+r_k)} + \cdots + \frac{1}{s_k(1+r_k)^{J/2}}\right]$$

$$\leq \frac{2m_k}{s_k\sqrt{2r_k \log(1+r_k)}}\sum_{j=0}^{\infty}\left(\frac{1}{1+r_k}\right)^{j/2} = \frac{2m_k}{s_k\sqrt{2r_k \log(1+r_k)}}\frac{\sqrt{1+r_k}}{\sqrt{1+r_k}-1}$$

$$\leq \left(\frac{m_k}{s_k}\right)\frac{4\sqrt{2}}{\sqrt{r_k \log(1+r_k)}}\frac{1}{r_k} \leq \frac{8m_k}{s_k r_k^2} \leq k. \qquad \square$$

COROLLARY 1. *Let $h_k$ be as defined in the proof of Theorem* 3. *Then*

$$\int f_0(y)\left(\log\frac{f_0(y)}{h_k(y)}\right)^2 dy \leq \omega_k \log V_k,$$

*where*

$$V_k = \left(1 + a_k\sqrt{S^2 + a_k^2}\right)(1 + r_k)^2(1 + \delta_k).$$

As noted in Wong and Shen (1995), it is sometimes useful to have the saturation rate in a distance that is stronger than Kullback–Leibler. Corollary 2 records the $\chi^2$ saturation rate. Recall that the $\chi^2$ distance is defined by

$$\chi^2(f, g) = \int \frac{(f - g)^2}{g}.$$

LEMMA 7. *If $f/g < V$ then*

$$\chi^2(f, g) \le 2(1 + V^{1/2})^2 d_{TV}(f, g).$$

PROOF.

$$\chi^2(f, g) = \int \frac{(f - g)^2}{g} = \int (\sqrt{f} - \sqrt{g})^2 \left(1 + \sqrt{\frac{f}{g}}\right)^2$$

$$\le (1 + V^{1/2})^2 d_H^2(f, g) \le 2(1 + V^{1/2})^2 d_{TV}(f, g). \qquad \square$$

COROLLARY 2. *Under the conditions of Theorem* 3,

$$\rho_k \equiv \inf_{f \in \mathcal{M}_k} \chi^2(f_0, f) \le 2(1 + V_k^{1/2})^2 \left\{2\delta_k\left[2 + \log(1 + a_k\sqrt{S^2 + a_k^2})\right] + 2r_k\right\},$$

*where*

$$V_k = \left(1 + a_k\sqrt{S^2 + a_k^2}\right)(1 + r_k)^2$$

*and $a_k$ is as defined in* (7).

PROOF. Define $f_k$ and $h_k$ as in Theorem 4. Then from the proof of Theorem 4 we see that

$$\frac{f_0}{h_k} = \frac{f_0}{f_k}\frac{f_k}{h_k} \le V_k.$$

Then, from Lemma 7 and Theorem 4,

$$\rho_k \le \chi^2(f_0, h_k) \le 2(1 + V_k^{1/2})^2 d_{TV}(f_0, h_k) \le 2(1 + V_k^{1/2})^2 D(f_0, h_k)$$

$$\le 2(1 + V_k^{1/2})^2 \left\{2\delta_k\left[2 + \log(1 + a_k\sqrt{S^2 + a_k^2})\right] + 2r_k\right\}. \qquad \square$$

**4. Rate of convergence.** Here we combine the results of the previous sections to compute the rate.

THEOREM 5. *Let $\alpha, \beta > 0$ and $\alpha_0, \beta_0 \geq 0$ be such that $2\alpha + \beta \leq 1$ and $2\alpha_0 + \beta_0 \geq 1$. Let $k_n = n^\beta/(\log n)^{\beta_0}$. Define $\omega_k$ as in (8) from Theorem 4 and let*

$$\text{(12)} \qquad t_n = \max\left\{ \frac{(\log n)^{\alpha_0}}{n^\alpha}, \frac{2}{c_1^{1/2}}\omega_{k_n}^{1/2} \right\}.$$

*Define $r_k$ as in Theorem 4. Then,*

$$P^*(d_H(f_0, \hat{f}_n) > t_n)$$

$$\text{(13)} \qquad \leq 4e^{-c_2 n t_n^2} + \frac{4\log\left((1 + a_{k_n}\sqrt{S^2 + a_{k_n}})(1 + r_{k_n})^2(1 + \delta_k)\right)}{c_1 n t_n^2}.$$

PROOF.   This proof parallels the argument in Theorem 3 of Wong and Shen (1995). Let $h_k \in \mathscr{M}_k$ be the density defined in the proof of Theorem 4. Let $B_n = \{f \in \mathscr{M}_{k_n}; d_H(f_0, f) > t_n\}$. Then,

$$P^*(d_H(f_0, \hat{f}_n) > t_n) \leq P^*\left( \sup_{f \in B_n^c} \prod_{i=1}^{n} \frac{f(Y_i)}{h_{k_n}(Y_i)} \geq e^{-c_1 n t_n^2/2} \right)$$

$$\leq P^*\left( \sup_{f \in B_n^c} \prod_{i=1}^{n} \frac{f(Y_i)}{f_0(Y_i)} \geq e^{-c_1 n t_n^2} \right)$$

$$+ P\left( \prod_{i=1}^{n} \frac{f_0(Y_i)}{h_{k_n}(Y_i)} \geq e^{-c_1 n t_n^2/2} \right) = P_1 + P_2.$$

Now, $P_1 \leq 4e^{-c_2 n t_n^2}$ by Theorem 3. We bound $P_2$ using Chebyshev's inequality, Theorem 4 and Corollary 1. Specifically, let $D_n = D(f_0, h_{k_n})$ and let

$$\gamma_n = \int f_0(y)\left( \log \frac{f_0(y)}{h_{k_n}(y)} \right)^2 dy.$$

Note that $D_n \leq \omega_{k_n} \leq (c_1/4)t_n^2$. Define $V_k$ as in Corollary 1. Then,

$$P\left( \prod_{i=1}^{n} \frac{f_0(Y_i)}{h_{k_n}(Y_i)} \geq e^{c_1 n t_n^2/2} \right) = P\left( \sum_{i=1}^{n} \log \frac{f_0(Y_i)}{h_{k_n}(Y_i)} \geq c_1 n t_n^2/2 \right)$$

$$= P\left( \sum_{i=1}^{n}\left( \log \frac{f_0(Y_i)}{h_{k_n}(Y_i)} - D_n \right) \geq n[c_1 t_n^2/2 - D_n] \right)$$

$$\leq \frac{n \operatorname{Var}\left( \log(f_0(Y)/h_{k_n}(Y)) \right)}{n^2(c_1 t_n^2/2 - D_n)^2} \leq \frac{16}{c_1^2 n}\frac{\gamma_n}{t^4}$$

$$\leq \frac{16}{c_1^2 n}\frac{\omega_{k_n}\log V_{k_n}}{t_n^4} \leq \frac{4}{c_1}\frac{\log V k_n}{n t_n^2}. \qquad \square$$

Now we consider some special cases.

4.1. *Compact support.*    An important special case studied in Roeder (1992) is when the mixing measure $P$ has compact support. Thus, suppose that $P(R) = 1$ where $R = \{(\mu, \sigma); s < \sigma < S, -m < \mu < m\}$ and $s, S$ and $m$ are positive constants. This class is still fully nonparametric but the conditions rule out infinite spikes and constrain the density $f_0$ to have thin tails.

First suppose that we take $m_k \to \infty$ and $s_k \downarrow 0$ in such a way that $m_k/s_k = k^\eta$ for some $\eta \in (0, 1)$. By Theorem 4, $\delta_{k_n} = 0$ for large $n$ so that $\omega_{k_n} \sim r_{k_n}$. Hence,

$$t_n \asymp \max \left\{ \frac{(\log n)^{\alpha_0}}{n^\alpha}, r_{k_n}^{1/2} \right\}$$

$$\asymp \max \left\{ \frac{(\log n)^{\alpha_0}}{n^\alpha}, \frac{(\log n)^{\beta_0(1-\eta)/4}}{\eta^{\beta(1-\eta)/4}} \right\}.$$

The expression is minimized by taking $\beta = 1 - 2\alpha$, $\beta_0 = 1 - 2\alpha_0$ and $\alpha = \alpha_0 = (1/2)(1 - \eta)/(3 - \eta)$ giving the rate

$$t_n \asymp \left( \frac{\log n}{n} \right)^{((1-\eta)/(3-\eta))/2}$$

which can be made arbitrarily close to $(\log n/n)^{1/6}$. We suspect that the $\log n$ can be eliminated by replacing the bracketing entropy with local entropy; see the comment after Theorem 1 of Wong and Shen (1995). In Section 5.1 we show how this rate can be improved to $(\log n/n)^{1/4}$.

Now consider choosing $m_k$ and $s_k$ so that $m_k/s_k = (\log k)^\eta$ for $\eta > 0$. Then, an analysis like that above yields the rate

$$t_n \asymp \max \left\{ \frac{(\log n)^{\alpha_0}}{n^\alpha}, \frac{(\log n)^{\beta_0/4}}{n^{\beta/4}} (\beta \log n - \beta_0 \log \log n)^{\eta/4} \right\}$$

$$\asymp \left( \frac{\log n}{n} \right)^{1/6} (\log n)^{\eta/6},$$

where the best rate is obtained by taking $\alpha = 1/6$, $\beta = 1 - 2\alpha$, $\alpha_0 = (1 + \eta)/6$ and $\beta_0 = 1 - 2\alpha_0$. In summary, choosing $k_n \sim n^{2/3}/(\log n)^{1/3}$ yields the rate $(\log n)^{(1+\eta)/6}/n^{1/6}$.

Now we consider data based choices of $m_k$ and $s_k$. A reasonable restriction on $m_{k_n}$ is $\hat{m}_n = \max |X_i|$ and it is easy to show that eventually $[-m, m] \subset [-\hat{m}_n, \hat{m}_n]$. From Roeder (1992), $\hat{m}_n = O(\sqrt{\log n})$ a.s. Next we estimate $s$. For this, we let $\hat{s}_n$ be the strongly consistent estimate from Theorem 4.2 of Roeder (1992), denoted by $\hat{h}_n$ in her paper. Her model is slightly different from ours but it can be shown her estimate of $s$ is still consistent in our setting. It follows that $\hat{m}_n/\hat{s}_n = O(\sqrt{\log n})$ almost surely for all large $n$. This corresponds to the above analysis with $\eta = 1/2$ giving a rate $(\log n)^{1/4}/n^{1/6}$.

4.2. *Noncompact support.*   To see how the tails of $P$ affect the rate, consider the simplified case where $s = S = 1$, say, so that $f_0(x) = \int \phi(x; \mu, 1) \, dP(\mu)$. Let $P$ have density $p$ with respect to Lebesgue measure. Suppose that $p$ is such that either (1) $p(\mu) \propto e^{-\lambda\mu}$ or (2) $p(\mu) \propto |\mu|^{-\lambda}$ for $\lambda > 0$ and $|\mu|$ large. In both these cases, the $a_k$ defined in Theorem 4 is infinite, which precludes us from finding a rate. Indeed, we conjecture that the estimate might even be inconsistent. This is not too surprising. To our knowledge, all sieve-based maximum likelihood estimates assume either a compactly supported density or a thin-tailed density. Similar, problems occur with kernel density estimates under Kullback–Leibler loss [Hall (1987)]. A remedy for this problem is discussed in the next section.

**5. Improving the rate.**   As we have seen, when the mixing $P$ does not have compact support, the rate of convergence is heavily affected by tail behavior of $P$. The sensitivity to the tails can be mitigated as follows. Let $\psi_0(x)$ be the density of a $t$-distribution with 1 degree of freedom and scale parameter $\widetilde{S} = 4\sqrt{\pi/2}$.

REMARK.   In general, we can take $\psi(x; \lambda, \mu, \sigma)$ to be the density of a $t$-distribution with $1/\lambda$ degrees of freedom centered on $\mu$ with scale, parameter $\sigma$. But this extra freedom does not appear to be needed.

Define a new sieve $\widetilde{\mathscr{M}}_{k_n}$ by

$$
(14) \quad \widetilde{\mathscr{M}}_k = \Bigg\{ f(\cdot) = p_0 \psi_0(x) + \sum_{j=1}^{k} p_j \phi(\cdot; \mu_j, \sigma_j); \\
|\mu_j| \leq m_k \text{ and } s_k \leq \sigma_j \leq S, \ j = 1, \dots, k \Bigg\}.
$$

THEOREM 6.   *Assume that the conditions of Theorem* 4 *hold except instead of $a_k$ we define $\tilde{a}_k = E(\mu^2/\sigma \mid R_k^c)$. Then,*

$$
(15) \quad \inf_{f \in \mathscr{M}_k} D(f_0, f) \leq 2\delta_k [2 + \log(1 + \tilde{a}_k)] + 2r_k \equiv \tilde{\omega}_k.
$$

PROOF.   For any $\mu$ and $\sigma > 0$,

$$
(16) \quad \frac{\phi_{\mu, \sigma}(x)}{\psi_0(x)} = \widetilde{S}\sqrt{\frac{\pi}{2}} \frac{1}{\sigma} \exp\left( -\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2 \right)[1 + (x - \mu + \mu)^2/\widetilde{S}^2].
$$

Note that $e^{-au^2}(u + b)^2$ for $a > 0$ and any $b$ attains a maximum that is less than or equal to $(b + 1/\sqrt{a})^2$. Hence,

$$
(17) \quad \frac{\phi_{\mu, \sigma}(x)}{\psi_0(x)} \leq \frac{1}{\sigma} 2\pi [1 + (\mu + \sigma\sqrt{2})^2/\widetilde{S}^2]
$$

$$
(18) \quad \leq \frac{\mu^2}{\sigma} \quad \text{for } |\mu| \geq 4 \text{ and } \sigma \leq 1.
$$

Now, if we replace $\phi_0$ by $\psi_0$ everywhere it appears in the proof of Theorem 4, we have the following. First,

$$(19) \qquad \frac{f_0}{f_k} \leq 1 + \frac{1}{\delta_k} \int_{R_k^c} \frac{\phi}{\psi_0}\, dP$$

$$(20) \qquad \leq 1 + \frac{1}{\delta_k} \int_{R_k^c} \frac{\mu^2}{\sigma}\, dP$$

$$(21) \qquad = 1 + \tilde{a}_k.$$

By Lemma 6, it follows that $D(f_0, f_k) \leq 2\delta_k[2 + \log(1 + \tilde{a}_k)]$. This bounds the first term in (9). The argument bounding the second term proceeds exactly as before with $\psi_0$ in place of $\phi_0$. This proves the theorem. □

Note that changing $\phi_0$ to $\psi_0$ in the sieve does not increase the bracketing entropy of the sieve, since $\psi_0$ is a fixed function. In other words, $\{\psi_0\}$ is a set of functions with bracketing number 1 (for all $\varepsilon$) and hence contributes a factor of 1 to the product in (6). We have immediately the following analogue of Theorem 5.

THEOREM 7.    *Let $\alpha, \alpha_0, \beta, \beta_0 > 0$ be such that $2\alpha + \beta \leq 1$ and $2\alpha_0 + \beta_0 \geq 1$. Let $k_n = n^\beta/(\log n)^{\beta_0}$. Define $\tilde{\omega}_k$ as in (15) from Theorem 4 and let*

$$(22) \qquad t_n = \max\left\{ \frac{(\log n)^{\alpha_0}}{n^\alpha}, \frac{2}{c_1^{1/2}} \tilde{\omega}_{k_n}^{1/2} \right\}.$$

*Define $r_k$ as in Theorem 4. Then,*

$$(23) \qquad P^*(d_H(f_0, \hat{f}_n) > t_n) \leq 4e^{-c_2 n t_n^2} + \frac{4\log((1 + \tilde{a}_{k_n})(1 + r_{k_n})^2)}{c_1 n t_n^2}.$$

PROOF.    Substitute $\widetilde{V}_k = (1 + \tilde{a}_k)(1 + r_k)^2$ for $V_k$ in the proof of Theorem 5 and use the $g_k$ as modified in Theorem 6. The calculation then proceeds exactly as before. □

5.1. *Compact support revisited.*    Now we show that, in the adjusted sieve, if $P$ has compact support, then the bound on $\inf_{f \in \mathscr{M}_k} D(f_0, f)$ can be improved. This leads to an improved rate of convergence.

THEOREM 8.    *Suppose that there exists $0 < s < S < \infty$ and $0 < m < \infty$ such that $P(\{(\mu, \sigma); -m \leq \mu \leq m, s < \sigma < S\}) = 1$. Then*

$$D_k \equiv \inf_{f \in \mathscr{M}_k} D(f_0, f) = O\left(\frac{\log k}{k}\right).$$

PROOF.    Let $f_0(x) = \int \phi(x; \mu, \sigma)\, dP(\mu, \sigma)$. Define a probability measure $Q$ on the real line by

$$Q((-\infty, t]) = \int_{-\infty}^t P\left(\mu + \sqrt{\sigma^2 - s^2}\, z \leq t\right) \phi(z)\, dz,$$

where $\phi(z) \equiv \phi(z; 0, 1)$. We claim that $f_0(x) = \int \phi(x; t, s) \, dQ(t)$. To see this, let $(\mu, \sigma)$, $Z_1, Z_2, Z_3$ be independent with $(\mu, \sigma) \sim P$ and $Z_1, Z_2, Z_3 \sim N(0, 1)$. Then,

$$X \overset{d}{=} \mu + \sigma Z_1 \overset{d}{=} \mu + \sqrt{\sigma^2 - s^2} Z_2 + s Z_3 \overset{d}{=} T + s Z_3,$$

where $T \sim Q$.

Set $c_k = m + \Delta_k$ where $\Delta_k^2 = 2b^2(\log k - \log \log k)$ and $b^2 = S^2 - s^2$. Let $f_k(x) = \int_{-c_k}^{c_k} \phi(x; t, s) \, dQ(t) + \delta_k \psi_0(x)$. Arguing as in Theorem 6,

$$D(f_0, f_k) \le 2\delta_k[2 + \log(1 + \tilde{a}_k)],$$

where $\delta_k = Q(|T| > c_k)$ and $\tilde{a}_k = s^{-1} E_Q(T^2 | |T| > c_k)$. Let $Z$ denote a standard normal random variable. Now,

$$
\begin{aligned}
Q(T > c_k) &= \Pr(\mu + \sqrt{\sigma^2 - s^2} Z > c_k) \\
&\le \Pr(\mu + \sqrt{\sigma^2 - s^2}|Z| > c_k) \\
&\le \Pr(m + \sqrt{S^2 - s^2}|Z| > m + \Delta_k) = \Pr\left(|Z| > \frac{\Delta_k}{b}\right) \\
&\le \frac{2b}{\sqrt{2\pi}\Delta_k} \exp\left\{-\frac{1}{2}\frac{\Delta_k^2}{b^2}\right\}.
\end{aligned}
$$

Hence,

$$\delta_k \le \frac{4b}{\sqrt{2\pi}\Delta_k} \exp\left\{-\frac{1}{2}\frac{\Delta_k^2}{b^2}\right\} = \frac{2}{\sqrt{\pi}} \frac{1}{(\log k - \log \log k)} \frac{\log k}{k} \le \frac{\log k}{k}$$

for large $k$.

By a similar argument,

$$E(T^2 \mid |T| > c_k) \le \frac{4b\Delta_k}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}\frac{\Delta_k^2}{b^2}\right\} \le \frac{4b^2 \log k}{\sqrt{\pi}k}.$$

So, for large $k$,

$$D(f_0, f_k) \le 2\delta_k[2 + \log(1 + \tilde{a}_k)] \le \frac{6\log k}{k}.$$

Let $B_k = 2(m + \Delta_k)/k$ and define $A_1 = [-c_k, -c_k + B_k)$, $A_2 = [-c_k + B_k, -c_k + 2B_k), \ldots, A_n = [c_k - B_k, c_k]$. Let $t_j$ be a point in $A_j$ and define $g_k(x) = \sum_j p_j \phi(x; t_j, s_k) + \delta_k \psi_0(x)$ where $p_j = \int_{A_j} dQ(t)$ and

$$s_k^2 = s^2 \left(1 + \frac{\log k}{k}\right).$$

Let $\phi_j(x) \equiv \phi(x; t_j, s)$. Then,

$$\int_{A_j} \phi(x; t, s) \, dQ(t) = \phi_j \int_{A_j} \frac{\phi(x; t, s)}{\phi_j} \, dQ(t)$$

$$\leq \phi_j \frac{s_k}{s} \int_{A_j} \exp\left\{ \frac{1}{2} \frac{(t - t_j)^2}{s_k^2 - s^2} \right\} dQ(t)$$

$$\leq p_j \phi_j \frac{s_k}{s} \exp\left\{ \frac{1}{2} \frac{B_k^2}{s_k^2 - s^2} \right\}$$

$$= p_j \phi_j \left\{ 1 + \frac{\log k}{k} \right\}^{1/2} \exp\left\{ \frac{1}{2} \frac{4(m + \Delta_k)^2 k}{k^2 s^2 \log k} \right\}$$

$$\leq p_j \phi_j \left\{ 1 + \frac{\log k}{k} \right\}^{1/2} \left( 1 + \frac{4(m + \Delta_k)^2 k}{k^2 s^2 \log k} \right)$$

for large $k$ since $e^x \leq 1 + 2x$ for $0 < x < 2$. For large $k$, we thus have that

$$\int_{A_j} \phi(x; t, s) \, dQ(t) \leq p_j \phi_j \left( 1 + \frac{\log k}{k} \right)^{3/2}.$$

Hence,

$$\frac{f_k}{g_k} \leq \frac{\sum_j \int_{A_j} \phi(x; t, s) \, dQ(t)}{\sum_j p_j \phi_j} \leq \left( 1 + \frac{\log k}{k} \right)^{3/2}.$$

Thus,

$$D(f_0, g_k) \leq D(f_0, f_k) + \sup_x \log \frac{f_k}{g_k} \leq \frac{15}{2} \frac{\log k}{k}$$

for large $k$. Finally, we must check that $g_k \in \mathcal{M}_k$. For this to be true, it suffices that the number of elements $N$ in the partition $A_1, \ldots, A_N$ is less than or equal to $k$. But $NB_k = \text{length}(-c_k, c_k) = 2(m + \Delta_k)$. So, $N = 2(m + \Delta_k)/B_k = k$ from the definition of $B_k$.   $\square$

Combining this result with Theorem 5 leads to the following.

COROLLARY 3.   *Assume the conditions in Theorem 8 above. Then choosing* $k_n \asymp \sqrt{n/\log n}$ *yields the rate of convergence* $\varepsilon_n \sim (\log n/n)^{1/4}$.

5.2. *Noncompact support case revisited.*   The revisited sieve allows us to compute a rate when $P$ has noncompact support. Consider again the simplified case where $s = S = 1$, so that $f_0(x) = \int \phi(x; \mu, 1) \, dP(\mu)$. Let $P$ have density $p$ with respect to Lebesgue measure such that $p$ has regularly varying tails, that is, either (1) $p(\mu) \propto e^{-\lambda|\mu|}$ or (2) $p(\mu) \propto |\mu|^{-\lambda}$, for $\lambda > 0$ and $|\mu|$ large.

In case (1),

(24)                      $\delta_k \propto e^{-\lambda m_k}$    and    $\tilde{a}_k \propto [1 + (\lambda m_k + 1)^2]$,

for any $\lambda > 0$. It follows that

$$(25) \qquad \tilde{\omega}_k = ce^{-\lambda m_k}\left\{2 + \log(1 + [1 + c'(\lambda m_k + 1)^2])\right\} + r_k$$

$$(26) \qquad \sim c''e^{-\lambda m_k}\log(m_k) + r_k.$$

for some constants $c, c', c'' > 0$ and where the second statement holds for large enough $m_k$. If $m_k = k^\eta$ for $0 < \eta < 1$, then

$$(27) \qquad \frac{e^{-\lambda m_k}\log(m_k)}{r_k} \propto \exp\left[\frac{1}{2}(1 - \eta)\ln k - \lambda k^\eta\right]\log(k) \to 0$$

as $k \to \infty$. hence, $\tilde{\omega}_k \asymp r_k$ and we recover the rates from the compact case. On the other hand, if $m_k = \log(k)$, then the size of $\lambda$ determines the dominant term. Specifically,

$$(28) \qquad \frac{e^{-\lambda m_k}\log(m_k)}{r_k} \propto k^{1/2-\lambda}\frac{\log\log k}{\sqrt{\log k}}.$$

If $\lambda \geq 1$, $r_k$ again dominates and we recover the rate from the compact case. If on the other hand $\lambda < 1$, then $\tilde{\omega}_k \propto k^{-\lambda}\log\log k$. Hence,

$$(29) \qquad \tilde{\omega}_{k_n}^{1/2} \propto n^{-\beta\lambda/2}(\log n)^{\beta_0\lambda/2}\sqrt{\log(\beta\log n - \beta_0\log\log n)}.$$

Choosing exponents to calibrate equation (22) gives us that

$$(30) \qquad t_n \asymp \left(\frac{\log n}{n}\right)^{(\lambda/(1+\lambda))/2}\sqrt{\log\log n}.$$

Similarly, in case (2),

$$(31) \qquad \delta_k \propto m_k^{1-\lambda} \quad \text{and} \quad \tilde{a}_k \propto m_k^2,$$

for $\lambda > 3$. If $\lambda \leq 3$, then $\tilde{a}_k$ is infinite, and we cannot compute a rate. It follows that when $\lambda > 3$,

$$(32) \qquad \tilde{\omega}_k = cm_k^{-(\lambda-1)}\left[2 + \log\left(1 + c'm_k^2\right)\right] + r_k$$

$$(33) \qquad \sim c''m_k^{-(\lambda-1)}\log(m_k) + r_k,$$

for some constants $c, c', c'' > 0$.

If $m_k = k^\eta$ for $0 < \eta < 1$, then $r_k$ dominates whenever $\lambda > ((1 + \eta)/\eta)$ and we recover the rates in the compact case. This happens whenever $\eta > 1/2$ because $\lambda > 3$. If $\lambda \leq ((1 + \eta)/\eta)$, then

$$(34) \qquad \tilde{\omega}_{k_n} \asymp k_n^{-\eta(\lambda-1)}\log k_n$$

$$(35) \qquad = (\log n)^{\beta_0\eta(\lambda-1)}n^{-\beta\eta(\lambda-1)}(\beta\log n - \beta_0\log\log n).$$

Calibrating exponents as above, the best rate is obtained with $\beta_0$ and $\beta = 1/[1 + \eta(\lambda - 1)]$. This yields

$$
(36) \qquad
\begin{aligned}
t_n &\asymp \left(\frac{\log n}{n}\right)^{(\eta(\lambda-1)/(1+\eta(\lambda-1)))/2} \log n^{(1/(1+\eta(\lambda-1)))/2} \\
&= n^{-(\eta(\lambda-1)/(1+\eta(\lambda-1)))/2}\sqrt{\log n}.
\end{aligned}
$$

Since $\lambda \le (1 + \eta)/\eta$ the exponent in the rate is no faster than $1/4$.

If $m_k = \log k$, then $m_{k_n}^{-(\lambda-1)} = (\beta \log n - \beta_0 \log\log n)^{-(\lambda-1)} \log(\beta \log n - \beta_0 \log\log n)$ while $r_{k_n} \propto (\beta \log n - \beta_0 \log\log n)^{1/2} n^{-1/2}$, so the first term in $\tilde{\omega}_{k_n}$ dominates for large enough $k_n$. Taking $\beta_0 = 0$, $\alpha_0 = 1/2$ and any $\alpha, \beta > 0$ satisfying the conditions of the theorem, we obtain

$$
(37) \qquad t_n \asymp (\log n)^{-(\lambda-1)/2}\sqrt{\log\log n},
$$

where the exponent is negative because $\lambda > 3$.

**6. Concluding remarks.**   We have established an upper bound on the rate of convergence for this mixture of Gaussian sieves. Our results suggest there is value in including long-tailed components in the sieve. The results are also interesting because the entropy calculations and saturation rate are nonstandard. We hope that these calculations will be useful for others working in the area of mixture asymptotics.

Finally, we mention three outstanding problems that form the subject of our current work. First, there is the question of whether the rate we have obtained is also a lower bound. (*Authors' note*: see the remark at the end of the introduction for recent developments on this question.) Second there is the problem of choosing the number of components $k_n$ from the data. We find the current methods for computing rates when the sieve index is chosen from the data—as in Barron and Yang (1995) for example—do not directly apply to finite mixtures. Third, we believe that some log terms in the rates can be eliminated by using local entropy instead of entropy. Again, for mixtures, calculating the local entropy appears to be nontrivial. We hope to report on these issues in a future paper.

## REFERENCES

BANFIELD, J. and RAFTERY, A. (1993). Model-based Gaussian and non-Gaussian clustering. *Biometrics* **49** 803–821.

BARRON, A. and YANG, Y. (1995). An asymptotic property of model selection criteria. Technical report, Dept. Statistics, Yale Univ.

CHEN, J. (1995). Optimal rate of convergence for finite mixtures models. *Ann. Statist.* **23** 221–233.

ESCOBAR, M. and WEST, M. (1995). Bayesian density estimation and inference using mixtures. *J. Amer. Statist. Assoc.* **90** 577–588.

GEMEN, S. and HWANG, C. (1982). Nonparametric maximum likelihood estimation by the method of sieves. *Ann. Statist.* **10** 401–414.

GHOSAL, S. and VAN DER VAART, A. (2000). Rates of convergence for Bayes and maximum likelihood estimation for mixtures of normal densities. Unpublished manuscript.

GRENANDER, U. (1981). *Abstract Inference*. Wiley, New York.

HALL, P. (1987). On Kullback–Leibler loss and density estimation. *Ann. Statist.* **15** 1491–1519.

LI, J. (1999). Estimation of mixtures models. Ph.D. dissertation, Dept. Statistics. Yale Univ.

LI, J. and BARRON, A. (1999). Mixture density estimation. Preprint.

LINDSAY, B. (1995). *Mixture Models: Theory, Geometry and Applications*. IMS, Hayward, CA.

MCLACHLAN, G. and BASFORD, K. (1988). *Mixture Models: Inference and Applications to Clustering*. Dekker, New York.

PRIEBE, C. (1994). Adaptive mixtures. *J. Amer. Statist. Assoc.* **89** 796–806.

ROBERT, C. (1996). Mixtures of distributions: inference and estimation. In *Markov Chain Monte Carlo in Practice* (W. Gilks, S. Richardson, D. Spiegelhalter, eds.) 441–464. Chapman and Hall, London.

ROEDER, K. and WASSERMAN, L. (1997). Practical Bayesian density estimation using mixtures of normals. *J. Amer. Statist. Assoc.* **92** 894–902.

ROEDER, K. (1992). Semiparametric estimation of normal mixture densities. *Ann. Statist.* **20** 929–943.

TONG, B., and VIELE, K. (1998). Mixtures of normal linear regressions. Technical report, Univ. Kentucky.

VAN DE GEER, S. (1996). Rates of convergence for the maximum likelihood estimator in mixture models. *Nonparametric Statist.* **6** 293–310.

VAN DER VAART, A. and WELLNER, J. (1996). *Weak Convergence and Empirical Processes*. Springer, New York.

WONG, W. and SHEN, X. (1995). Probability inequalities for likelihood ratios and convergence rates of sieve MLEs. *Ann. Statist.* **23** 339–362.

DEPARTMENT OF STATISTICS
232 BAKER HALL
CARNEGIE MELLON UNIVERSITY
PITTSBURGH, PENNSYLVANIA 15213
E-MAIL: genovese@stat.cmu.edu