

On Selection Rules for Finite Mixtures of  
Distributions\*

by

Shanti S. Gupta, Purdue University

and

Wen-Tao Huang, Academia Sinica, Taipei

Department of Statistics  
Division of Mathematical Sciences  
Mimeograph Series #504

August 1977

\*This research was supported by the Office of Naval Research Contract N00014-75-C-0455 at Purdue University. Reproduction in whole or in part is permitted for any purpose of the United States Government.

On Selection Rules for Finite Mixtures of  
Distributions\*

by

Shanti S. Gupta, Purdue University  
and

Wen-Tao Huang, Academia Sinica, Taipei

1. Introduction and Summary

In some industrial processes, ordinarily observations arise from some distribution function  $F_1(x;\theta)$ , say, but occasionally, the process yields "outliers" which may follow some other distribution  $F_2(x;\varphi)$ . Accordingly, if outliers carry no obvious label, then the process produces observations according to a mixture distribution  $\alpha F_1(x;\theta) + (1-\alpha)F_2(x;\varphi)$  for some proportion  $\alpha$ ,  $0 < \alpha < 1$ . Also, in marine biology, one may be interested in studying certain characteristics of a fish. For this purpose, samples of fish are taken and the desired trait is measured for each fish. Since many characteristics vary according to the age of fish, the trait has a distinct distribution for each age group and the population has a mixture of distributions. On the other hand, mixtures of distributions occur in the compound decision problems as proposed by Robbins [8], in which mixing distributions correspond to some a priori distributions.

It happens that in many cases, an experimenter is faced with a problem of choosing one or more "desirable" processes (treatments) from among  $k$  given processes (treatments) which produce observations according to some

---

\*This research was supported by the Office of Naval Research contract N00014-75-C-0455 at Purdue University. Reproduction in whole or in part is permitted for any purpose of the United States Government.

mixture distributions. For his special purpose, the experimenter may need one or more processes which are associated with the largest (smallest) proportion in the mixtures of distributions. For instance, he may need the process which has the least proportion of occurrence of outliers.

The problem for the estimation of these proportions of mixtures is not easy. For example, when  $F_1(x; \theta)$  and  $F_2(x; \varphi)$  both are normal with a common variance  $\sigma^2$  and, with means  $\mu_1$  and  $\mu_2$ , respectively, if  $\mu_1$  and  $\mu_2$  are not well separated, i.e. when  $d \equiv |\mu_1 - \mu_2|/\sigma$  is small, it is almost impossible to classify the observations from the mixture distribution into two groups. To be more precise, let  $I(\alpha; F_1, F_2)$  denote the Fisher information for the estimation of  $\alpha$ . Hill [ 5 ] pointed out that for  $d$  small,  $I(\alpha; F_1, F_2) \approx d^2$ . Therefore, if  $d$  is in  $(1/8, 1/4)$ , then for the maximum likelihood estimate for  $\alpha$  with standard deviation 0.1, the sample size needed is large, as big as 6400. However, so far the classical efficient method for the estimation of  $\alpha$  is the maximum likelihood estimate. The usual moment estimate is inefficient. However, when the number of components of mixture increases, situations become more complicated even for the maximum likelihood estimates. This suggests that another approach should be considered. The so-called minimum distance method of Wolfowitz [12] seems reasonable. Large sample properties like consistency can be shown to hold. If the distance between two distribution functions is properly chosen, some other optimal properties may hold. And, if the rate of convergence is fair, then this approach should be right. In the problems of selection and ranking, some statistics are necessary such that based on these quantities, the criterion of priority of selection can be constructed. Though these statistics may not necessarily be good estimates

for the unknown parameters which are under consideration, most of them may do well for the selection problem. Accordingly, the minimum distance method seems natural to be one of the approaches to follow for the problems of selecting the largest (or smallest) proportion of certain component of mixture. The so-called least squares method will be applied in this paper for the selection problem.

In section 2 some notation are defined and the problem is formulated. A class of consistent selection procedures are defined in section 3 and some asymptotic optimal properties are shown.

## 2. Notation and Formulation of the Problem

The problem of identifiability should be mentioned, since the selection problem for the proportion of mixtures is related to the identifiability problem. This can be simply illustrated by an example. Let  $B(n;p)$  denote a binomial distribution with success probability  $p$ , then, it can be found some  $\alpha_1, \alpha_2, \beta_1, \beta_2$  and some  $p_1, p_2$  and  $p_3$  such that  $\alpha_1 \neq \beta_1$  and  $\alpha_1 B(n;p_1) + \alpha_2 B(n;p_2) + (1-\alpha_1-\alpha_2)B(n;p_3)$  and  $\beta_1 B(n;p_1) + \beta_2 B(n;p_2) + (1-\beta_1-\beta_2)B(n;p_3)$  represent the same mixture distribution if  $n < 5$ . In this example, it is impossible to identify and select. Necessary and sufficient conditions for identifiability of finite mixtures can be found in [11] and [13].

Let  $\mathfrak{F}$  denote the family of distributions such that the associated convex hull of  $\mathfrak{F}$  is identifiable. Many well-known families of distributions are included in  $\mathfrak{F}$ . For example (see [11], [13]),  $\mathfrak{F}$  can be family of  $p$ -variate normal distributions, product of  $n$  exponential distributions, binomial distributions with different integral parameters, translation

parameter family induced by a certain univariate cdf, union of the families of product of  $n$  exponential distributions and the  $p$ -variate normal distribution etc.

For convenience, for some prefixed integer  $m$ , we define

$$(2.1) \quad \langle 0,1 \rangle^m = \{(\alpha_1, \alpha_2, \dots, \alpha_m) : \alpha_i > 0, \sum_{i=1}^m \alpha_i = 1\} \quad (m \geq 2).$$

Let  $\lambda$  be a real-valued continuous function on  $\langle 0,1 \rangle^m$ . Let the functions  $F_1(x; \theta_1), F_2(x; \theta_2), \dots, F_m(x; \theta_m)$  be in  $\mathcal{F}$ , where  $\theta_i$  may be a parameter vector and  $F_i(x; \theta_i)$  and  $F_j(x; \theta_j)$  may have different parametric form, for instance,  $F_i(x; \theta_i)$  may be a normal distribution with location-scale parameter  $(\mu_i, \sigma_i^2)$  and  $F_j(x; \theta_j)$  may be an exponential distribution with location-scale parameter  $(\alpha_j, \beta_j)$ . For convenience, we denote

$$(2.2) \quad \underline{F} = (F_1(x; \theta_1), \dots, F_m(x; \theta_m))$$

and

$$(2.3) \quad \underline{\alpha}_i = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{im}).$$

A finite mixture distribution with  $m$  component is defined to be the inner product of certain  $\underline{\alpha} \in \langle 0,1 \rangle^m$  and  $\underline{F}$ , i.e.

$$(2.4) \quad G(x; \underline{\alpha}) = \underline{\alpha} \cdot \underline{F} \\ = \sum_{i=1}^m \alpha_i F_i(x; \theta_i)$$

Let  $\pi_1, \pi_2, \dots, \pi_k$  be  $k$  populations such that  $\pi_i$  has cdf  $G(x; \underline{\alpha}_i)$  (defined by (2.4)) for some unknown parameter  $\underline{\alpha}_i \in \langle 0,1 \rangle^m$ . Let  $X_{i1}, X_{i2}, \dots, X_{im}$  be  $n$

random observations from  $\pi_i$ ,  $i=1,2,\dots,k$ . Let  $G_{in}(x)$  denote the associated empirical distribution function. Let  $\lambda_{[1]}(\alpha) \leq \lambda_{[2]}(\alpha) \leq \dots \leq \lambda_{[k]}(\alpha)$  denote the order values of  $\lambda(\alpha_1), \lambda(\alpha_2), \dots, \lambda(\alpha_k)$ .

Based on  $n$  independent observations from each population, we are interested in selecting  $t$  ( $1 \leq t \leq k-1$ ) populations, say,  $\pi_{r_1}, \pi_{r_2}, \dots, \pi_{r_t}$  such that  $\lambda(\alpha_{r_1}), \lambda(\alpha_{r_2}), \dots, \lambda(\alpha_{r_t})$  are the  $t$  largest. We call these populations the  $t$  best.

We approach the problem by the indifference zone formulation. For convenience, we introduce the following notation.

For given  $\Delta$ , we define

$$(2.5) \quad \Omega(\lambda; \Delta) = \{(\alpha_1, \alpha_2, \dots, \alpha_k) : \alpha_i \in \langle 0, 1 \rangle^m, \lambda_{[k-t+1]}(\alpha) \geq \lambda_{[k-t]}(\alpha) + \Delta\}.$$

For specified  $F$  and  $\lambda$ , we consider our problem on the configuration  $\Omega(\lambda; \Delta)$  for given  $\Delta$  for the indifference zone approach. We also define

$$(2.6) \quad \Omega = \langle 0, 1 \rangle^m \times \langle 0, 1 \rangle^m \times \dots \times \langle 0, 1 \rangle^m. \quad (k \text{ copies})$$

Finally, we define, for given  $p$ ,  $0 \leq p \leq 1$

$$(2.7) \quad S(\alpha; H) = \int_{-\infty}^{\infty} (\alpha \cdot F - G_n(x))^2 dH(x)$$

where  $\alpha \cdot F$  is a mixture distribution for  $\alpha \in \langle 0, 1 \rangle^m$  and  $G_n(x)$  is the empirical distribution associated with some  $\alpha_0 \cdot F$  for unknown  $\alpha_0 \in \langle 0, 1 \rangle^m$ . And  $H(x)$  is a cdf. Hence,  $S(\alpha; H)$  is a function on  $\alpha \in \langle 0, 1 \rangle^m$ .

### 3. A class of consistent selection procedures

In this section, we consider the cases when  $F$  are continuous and discrete. In each case, we assume the component  $F_i(x_i; \theta_i)$  of  $F$  are completely known.

## (A) Continuous case

We assume the parametric form of each component  $F_i(x; \theta_i)$  of  $F$  is continuous in  $x$  for each  $\theta_i$  and continuous in  $\theta_i$  for each  $x$ .

For given  $n$  observations from a population with cdf  $G(x; \alpha_0) = \alpha_0 \cdot F$  for some unknown  $\alpha_0$  and a given cdf  $H(x)$ , a vector  $\hat{\alpha} \in \langle 0, 1 \rangle^m$  at which  $S(\alpha; H)$  attains its infimum seems a "good" estimate for the real  $\alpha_0$  in the sense of least squares method. It is to be noted that  $\hat{\alpha}$  is a statistic of  $n$  observations and also is a function of  $F$  and  $H$ . A good choice in some sense for the weight function  $H(x)$  is not easy. Bartlett and Macdonald [1] study some special case of  $m = 2$ . For  $m \geq 3$ , the situation is complicated.

Choi and Balgren [3] consider the case  $H(x) = G_{in}(x)$  and obtain some optimal properties like consistency and asymptotical normality. However, for the case of small samples, Macdonald [7] points out that, using  $H(x) = \alpha \cdot F$ , some Monte Carlo results show some improvement of the Choi and Bulgren's result. And, as a matter of fact, for  $H(x) = \alpha \cdot F$ ,  $S(\alpha; H)$  is the von Mises statistic for the goodness-of-fit. Let  $U_1$  and  $U_2$  denote two random observations from the population with cdf  $F(x)$  and  $V_1$  and  $V_2$  denote the random observations from a population with cdf  $G(x)$ . It is known that  $\Delta(F, G) \equiv P_r \{ U_1 \vee U_2 < V_1 \wedge V_2 \text{ or } V_1 \vee V_2 < U_1 \wedge U_2 \} = 1/3 + 1/2 \int_{-\infty}^{\infty} (F(x) - G(x))^2 d(\frac{F(x)+G(x)}{2})$  where  $a \vee b = \max(a, b)$ ,  $a \wedge b = \min(a, b)$  (Lehmann [6]). Note that  $\Delta(F, G) = 0$  if, and only if  $F \equiv G$ . Roughly speaking, taking  $F(x)$  to be  $\alpha \cdot F$  and  $G(x)$  to be  $G_m(x)$ , it is significant to consider  $H(x) = \frac{1}{2} (\alpha \cdot F + G_m(x))$  for our case. Accordingly, in general, we consider the case  $H(x) = p \alpha \cdot F + (1-p)G_m(x)$  for  $0 \leq p \leq 1$ . Note that

$p = 0$  yields the Choi and Bulgren's case and for  $p = 1$ , we get the Macdonald's case. For our notational convenience, henceforth, we define

$$(3.1) \quad S_i(\alpha; p) = \int_{-\infty}^{\infty} (\alpha \cdot F - G_{in}(x))^2 d(p \alpha \cdot F + (1-p)G_{in}(x))$$

which is obtained by taking  $H(x) = p \alpha \cdot F + (1-p)G_{in}(x)$  where  $G_{in}(x)$  is the empirical distribution associated with the  $n$  random observations from the population  $\pi_i$ . The existence of some  $\hat{\alpha}_i$  such that  $S_i(\alpha; p)$  attains the infimum can be shown by going through the analogous arguments as in [3]. Define  $\hat{\alpha}_i$  to be such that

$$(3.2) \quad S_i(\hat{\alpha}_i; p) = \inf_{\alpha \in \langle 0, 1 \rangle^m} S_i(\alpha; p).$$

For a given value of  $p$  ( $0 \leq p \leq 1$ ), we define a selection procedure  $R_p$  as follows.

Take  $n$  independent observations from each  $\pi_i$  and construct the empirical distribution  $G_{in}(x)$ . Compute  $\hat{\alpha}_i = \hat{\alpha}_i(X_{i1}, X_{i2}, \dots, X_{in})$  which is defined by (3.1) and (3.2). Let  $\lambda_{[1]}(\hat{\alpha}) \leq \lambda_{[2]}(\hat{\alpha}) \leq \dots \leq \lambda_{[k]}(\hat{\alpha})$  denote the ordered values of  $\lambda(\hat{\alpha}_1), \lambda(\hat{\alpha}_2), \dots, \lambda(\hat{\alpha}_k)$ .

$R_p$ : Select  $\pi_i$  if, and only if  $\lambda(\hat{\alpha}_i) \geq \lambda_{[k-t+1]}(\hat{\alpha})$

Use a mechanism when a tie occurs.

By a correct selection (CS) we mean a set of  $t$  populations associated with the  $t$  largest values of  $\lambda(\hat{\alpha}_1), \lambda(\hat{\alpha}_2), \dots, \lambda(\hat{\alpha}_k)$  is selected.

**Definition 3.1** A selection procedure  $R$  is consistent with respect to

$$(\mathfrak{F}, \lambda) \text{ if } \lim_{\Delta \rightarrow 0} \lim_{n \rightarrow \infty} \inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{CS|R\} = 1$$

**Definition 3.2** A selection procedure  $R$  is strongly asymptotically monotone with respect to  $(\mathfrak{F}, \lambda)$  if  $\lambda(\alpha_i) < \lambda(\alpha_j)$  and for any  $\epsilon > 0$  implies



$$\lim_{n \rightarrow \infty} \sup_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{ \pi_i \text{ is selected } | R \} - \epsilon < \lim_{n \rightarrow \infty} \sup_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{ \pi_j \text{ is selected } | R \}.$$

**Theorem 3.1** For any value of  $p$ ,  $0 \leq p \leq 1$ ,  $R_p$  is consistent and strongly asymptotically monotone with respect to  $(\mathcal{F}, \lambda)$ .

**Proof:** (a) We show that for any  $p$  ( $0 \leq p \leq 1$ ) and for each  $i$  ( $i=1,2,\dots,k$ ),  $\hat{\alpha}_{\sim i} \rightarrow \alpha_{\sim i}$  with probability one. We follow the arguments given in [3] with appropriate modifications. Now, by the Glivenko-Cantelli theorem, for  $\epsilon > 0$ ,  $\exists N(\epsilon)$  such that, whenever  $n \geq N(\epsilon)$ ,

$$P\{|p \alpha_{\sim i} \cdot F + (1-p)G_{in}(x) - G_{in}(x)| < \epsilon\} = P\{|p \alpha_{\sim i} \cdot F - G_{in}(x)| < \epsilon\} = 1.$$

Replacing  $dF_n(x)$  by  $d(p\alpha_{\sim i} \cdot F + (1-p)G_{in}(x))$  and following the same argument as given in the proof of Theorem 2 in [3], the result follows.

(b) Consistency of  $R_p$

Since  $\lambda$  is continuous it is true that  $\lambda(\hat{\alpha}_{\sim i}) \rightarrow \lambda(\alpha_{\sim i})$  WPI ( $i=1,2,\dots,k$ ). Now, by the Egoroff's theorem, for  $\epsilon > 0$  and  $\delta > 0$ , there exists  $N_i(\epsilon, \delta)$ ,  $A_i$  and  $B_i$  such that the sample space is decomposed to be  $A_i \cup B_i$  with  $B_i$  the complement of  $A_i$  and  $P(B_i) > 1 - \epsilon$  and on  $B_i$ ,  $|\lambda(\hat{\alpha}_{\sim i}) - \lambda(\alpha_{\sim i})| < \delta$  whenever  $n \geq N_i(\epsilon, \delta)$  uniformly in  $\alpha_{\sim i} \in (0,1)^m$ , i.e.,  $N(\epsilon, \delta)$  is independent of  $\alpha_{\sim i}$ . Note that  $\lambda(\hat{\alpha}_{\sim i})$  depends on  $n$ . Set  $N = N_1(\epsilon, \delta) + \dots + N_k(\epsilon, \delta)$  and set  $B = \bigcap_{i=1}^k B_i$ . Then,  $P(B) > 1 - \epsilon$ , and on  $B$ , whenever  $n \geq N$ ,  $\max_{1 \leq i \leq k} |\lambda(\hat{\alpha}_{\sim i}) - \lambda(\alpha_{\sim i})| < \delta$  uniformly for each  $(\alpha_{\sim 1}, \alpha_{\sim 2}, \dots, \alpha_{\sim k}) \in \Omega$  (defined by (2.6)). Now, for any given  $P^* \in (0,1)$ , and for given  $\Delta > 0$ , however small, choose  $\delta = \frac{\Delta}{3} > 0$  and  $\epsilon = 1 - P^*$ . Since on  $\Omega(\lambda; \Delta)$ ,  $\lambda_{[k-t+1]} - \lambda_{[k-t]} \geq \Delta = 3\delta$ . Hence we conclude that

$$P_{\alpha} \{ \lambda(\hat{\alpha}_{\sim r_i}) > \lambda_{[k-t]}(\alpha), i=1,2,\dots,t | \lambda(\alpha_{\sim r_i}) > \lambda_{[k-t]}(\alpha) \} > P^*$$

for every  $\alpha \in \Omega(\lambda; \Delta)$ . Hence, we have shown that for every  $\Delta > 0$ ,  $\lim_{n \rightarrow \infty}$

$\inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{CS | R_p\} = 1$ . Hence the consistency is shown.

(c) Suppose  $\lambda(\alpha_i) < \lambda(\alpha_j)$ .

(i) If  $\lambda(\alpha_i) \leq \lambda_{[k-t]}(\alpha)$  and  $\lambda(\alpha_j) \geq \lambda_{[k-t+1]}(\alpha)$ . Then, take  $P^* \geq 2/3$  and go through the arguments given in (b), we conclude that  $\inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{\pi_j$

is selected  $|R_p\} \geq \inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{CS | R_p\} \geq 2/3$  whenever  $n \geq N_0 = N_0(\Delta)$

for some  $N_0$ . On the other hand, for each  $n \geq N_0$ ,  $\{\pi_i$  is selected  $|R_p\}$

$\subset$  {Selection is not correct  $|R_p\}$ . Hence  $P_{\alpha} \{\pi_i$  is selected  $|R_p\} \leq 1 - P_{\alpha} \{CS | R_p\} \leq 1/3 \quad \forall \alpha \in \Omega(\lambda; \Delta)$ , i.e.

$$\sup_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{\pi_i \text{ is selected } |R_p\} \leq 1/3 \text{ for each } n \geq N_0.$$

(ii) Suppose both  $\lambda(\alpha_i)$  and  $\lambda(\alpha_j)$  are no larger than  $\lambda_{[k-t]}(\alpha)$ .

Then, for  $\epsilon > 0$  and by the arguments in (b), there exists a subset of sample space  $B$  and an integer  $N_0$  such that  $P\{B\} > 1 - \frac{\epsilon}{2}$  and for  $n \geq N_0$

and on  $B$ ,  $\max_{1 \leq i < k} \{|\alpha_i - \hat{\alpha}_i|\} < \frac{\Delta}{3}$ . Let  $E$  denote the event  $\{\pi_i$  is

selected  $|R_p\}$ . Then  $E = E \cap B + E \cap B^c$ . Hence,  $\sup_{\alpha} P_{\alpha} \{E\} \leq \sup_{\alpha} P_{\alpha} \{E \cap B\}$

$$+ \sup_{\alpha} P_{\alpha} \{E \cap B^c\} \leq \sup_{\alpha} P_{\alpha} \{E \cap B\} + \frac{\epsilon}{2}.$$

since  $P_{\alpha} \{E \cap B^c\} \leq P_{\alpha} \{B^c\} < \frac{\epsilon}{2} \quad \forall \alpha \in \Omega(\lambda; \Delta)$ . Noting that for any

$\alpha \in \Omega(\lambda; \Delta)$ ,  $P_{\alpha} \{E \cap B\} = 0$  since on  $B$ ,  $\hat{\alpha}_i < \alpha_{[k-t+1]} - \frac{\Delta}{3}$ .

(iii) If  $\lambda(\alpha_i)$  and  $\lambda(\alpha_j)$  both are no less than  $\lambda_{[k-t+1]}(\alpha)$ . The proof is analogous to the case of (ii).

The proof is thus complete.

**Remark 3.1.** If  $t_1, t_2, \dots, t_m$  are some positive integers such that each  $t_i$  is no larger than  $k - 1$ . Let  $\Omega(t_1, t_2, \dots, t_m) \equiv \{(\alpha_{\sim 1}, \alpha_{\sim 2}, \dots, \alpha_{\sim k}) : \alpha_{[k-t_i+1]}^{(i)} > \alpha_{[k-t_i+1]}^{(i)} \ i = 1, 2, \dots, m\}$  where  $\alpha_{[j]}^{(i)}$  denotes the  $j$ -th largest value of the  $i$ -th component of  $\alpha_{\sim 1}, \alpha_{\sim 2}, \dots, \alpha_{\sim k}$  and we denote  $\alpha_{\sim r} = (\alpha_{\sim r}^{(1)}, \alpha_{\sim r}^{(2)}, \dots, \alpha_{\sim r}^{(m)})$ . If for each  $i$  we are desired to select the  $t_i$  largest in the  $i$ -th component simultaneously, then, using the statistics  $\{\hat{\alpha}_{\sim 1}, \hat{\alpha}_{\sim 2}, \dots, \hat{\alpha}_{\sim k}\}$ , which are defined by (3.2), associated with the  $i$ -th component, we select these populations which have the  $t_i$  largest values in the  $i$ -th component of  $\{\alpha_1^{(i)}, \alpha_2^{(i)}, \dots, \alpha_k^{(i)}\}$  ( $i=1, 2, \dots, m$ ). It can be shown that the simultaneous selections are also consistent and strongly asymptotically monotone on the configuration  $\Omega(t_1, t_2, \dots, t_k)$ .

**Remark 3.2.** For  $m = 2$  and a given  $n$ , let  $\hat{\alpha}_i'$ ,  $\hat{\alpha}_i''$  and  $\hat{\alpha}_i$  denote respectively, the least square estimates associated with  $p = 0$ ,  $p = 1$  and some  $p(0 < p < 1)$ . Then, it can be obtained

$$\hat{\alpha}_i' = \frac{\Sigma(F_2 - F_1)(F_1 - \frac{i}{n})}{\Sigma(F_2 - F_1)^2}, \quad \hat{\alpha}_i'' = \hat{\alpha}_i' + \frac{1}{2n} \frac{\Sigma(F_2 - F_1)}{\Sigma(F_2 - F_1)^2}$$

and

$$\hat{\alpha}_i = \hat{\alpha}_i' + \frac{1-p}{2n} \frac{\Sigma(F_2 - F_1)}{\Sigma(F_2 - F_1)^2}$$

where

$$\Sigma(F_2 - F_1) \equiv \sum_{i=1}^n (F_2(X_{[i]}; \theta_2) - F_1(X_{[i]}; \theta_1)) \text{ and } X_{[1]} < X_{[2]} < \dots < X_{[n]}$$

are the order statistics from  $\pi_i$ . As a convention we take  $\hat{\alpha}_i' = 0$  if  $\hat{\alpha}_i' < 0$  and  $= 1$  if  $\hat{\alpha}_i' \geq 1$  and use the same convention for other two cases. It can be seen that  $\hat{\alpha}_i$  is always between  $\hat{\alpha}_i'$  and  $\hat{\alpha}_i''$  for all  $n$ . If  $F_1$  and

$F_2$  are "smooth" in some sense, we see that  $|\hat{\alpha}_i' - \hat{\alpha}_i''| = O(n^{-1+\epsilon})$  for  $\epsilon > 0$ .

**Definition 3.3** A selection procedure  $R$  is consistent of order  $O(A(\Delta))$  ( $o(A(\Delta))$ ) with respect to  $(\mathcal{F}, \lambda)$  if

$$\lim_{\substack{\Delta \rightarrow 0 \\ n=O(A(\Delta))}} \inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{CS|R\} = 1 \quad (\lim_{\substack{\Delta \rightarrow 0 \\ n=O(A(\Delta))}} \inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{CS|R\} = 1).$$

**Theorem 3.2** For given  $p$ ,  $0 \leq p \leq 1$ ,  $R_p$  is consistent of order  $O(\Delta^{\frac{-2}{1-2\delta}})$ ,  $0 < \delta < 1/2$ .

**Proof:** We note that, by the Glivenko-Cantelli theorem that

$\sup_x |G_i(x) - G_{in}(x) + o(1)| \rightarrow 0$  WPl as  $n \rightarrow \infty \quad \forall_i$ , where  $o(1)$  is independent of  $x$ . For any fixed  $i$  ( $1 \leq i \leq k$ ), let  $\dot{S}(\underline{\alpha}_i; p)$  denote the  $m-1$  equations for which each equation is differentiated with respect to  $\alpha_{ij}$ ,  $j = 1, 2, \dots, m-1$ , where  $\underline{\alpha}_i = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{im-1}, 1 - \sum_{j=1}^{m-1} \alpha_{ij})$ . Then, the first element of  $\dot{S}(\underline{\alpha}_i; p)$  for  $j=1$  becomes

$$\begin{aligned} & \frac{1}{n} \sum_{j=1}^n F_1(X_{i[j]}; \theta_1) \left\{ \sum_{r=1}^m \alpha_{ir} F_r(X_{i[j]}; \theta_r) \frac{j}{n} + \frac{1-p}{2n} \right\} \\ & \leq \sup_x |G_i(x) - G_{im}(x) + o(1)| \frac{1}{n} \sum_{j=1}^n F_1(X_{i[j]}; \theta_1). \end{aligned}$$

where  $X_{i[j]} \leq \dots \leq X_{i[n]}$  are order statistics from  $\pi_i$ . Follow the analogous arguments of the proof of Theorem 4 of [3], we conclude that  $|\hat{\alpha}_i - \underline{\alpha}_i| < O(n^{-\frac{1}{2}+\delta})$  for all but finite  $n$  with prob. 1 where  $0 \leq \delta < \frac{1}{2}$ . Now, if we take  $\Delta/2 = O(n^{-\frac{1}{2}+\delta})$  for large  $n$ , we see that  $n = O(\Delta^{\frac{-2}{1-2\delta}})$

and for this  $n$  it can be sure that the selection is correct with probability one as  $\Delta \rightarrow 0$ . The proof is thus complete.

Let  $\bar{\alpha}_{\sim i}$  denote the arithmetic mean of  $r$  independent estimates of  $\hat{\alpha}_{\sim i}$  where  $r$  is some integer. This means  $rn$  samples are drawn from each population. And for each subgroups of  $n$  samples, we obtain an estimate  $\hat{\alpha}_{\sim i}$  for the population  $\pi_i$ . If  $n$  is large,  $\lambda(\alpha_i) = \alpha_{i1}$ , and  $t = 1$ , we propose the following rule  $R'_p$ .

$R'_p$ : Select  $\pi_i$  if  $\bar{\alpha}_{\sim i1} \geq \bar{\alpha}_{\sim j1}$  for all  $j \neq i$ .  
where  $\bar{\alpha}_{\sim i1}$  is the first component of  $\bar{\alpha}_{\sim i}$ .

Theorem 3.3 If  $n$  is large,  $t = 1$ , and  $\lambda(\alpha_i) = \alpha_{i1}$ , the projection function, then we have

$$\inf_{\alpha \in \Omega(\lambda; \Delta)} P_{\alpha} \{CS | R'_p\} \geq \int_{-\infty}^{\infty} \prod_{j=2}^k \Phi(\delta_j z + \frac{\sqrt{r} \Delta}{\sigma[j]}) d\Phi(z)$$

where  $\Phi(x)$  denotes the standard normal distribution and

$$\sigma_j^2 = 2 \int_{-\infty < x < y < \infty} G_j(x) [1 - G_j(y)] dB_j(x) dB_j(y)$$

where

$$B_j(x) = F_1(x; \theta_1) G_j(x) - \int_{-\infty}^{\infty} F_1(x; \theta_1) dG_j(x)$$

for  $j=1, 2, \dots, k$ .

and  $\sigma[1] \leq \sigma[2] \leq \dots \leq \sigma[k]$ ,  $\delta_j = \alpha[1] / \sigma[j]$ .

Proof: It has been shown in [2] that  $\hat{\alpha}_{\sim i}$  is asymptotically normal and hence, the first component of  $\hat{\alpha}_{\sim i}$ , say,  $\hat{\alpha}_{\sim i1}$  is asymptotically normal with mean  $\alpha_{i1}$  and variance  $\sigma_i^2 = 2 \int_{-\infty < x < y < \infty} G_i(x) [1 - G_i(y)] dB_i(x) dB_i(y)$  where

$$B_i(x) = F_1(x; \theta_1) G_i(x) - \int_{-\infty}^{\infty} F_1(x; \theta_1) dG_i(x).$$

Hence, when  $n$  is large,  $t = 1$ , we have for  $\alpha \in \Omega(\lambda; \Delta)$

$$\begin{aligned} P_{\alpha} \{CS | R'_p\} &= P_{\alpha} \{ \bar{\alpha}_{k1} \geq \bar{\alpha}_{j1} \quad j=1, 2, \dots, k-1 \mid \alpha_{k1} = \max_{1 \leq j \leq k} \alpha_{j1} \} \\ &= P_{\alpha} \left\{ \frac{\sqrt{r}(\bar{\alpha}_{k1} - \alpha_{k1})}{\sigma_k} \geq \frac{\sqrt{r}(\bar{\alpha}_{j1} - \alpha_{j1})}{\sigma_j} \frac{\sigma_j}{\sigma_k} + \frac{\sqrt{r}(\alpha_{j1} - \alpha_{k1})}{\sigma_k} \right\} \\ &\geq P_{\alpha} \left\{ Z_k \geq Z_j \left( \frac{\sigma_j}{\sigma_k} \right) - \frac{\sqrt{r} \Delta}{\sigma_k} \quad j=1, 2, \dots, k \right\} \text{ (where } Z_1, Z_2, \dots, Z_k \text{ are} \\ &\quad \text{iid standard normal)} \\ &= \int_{-\infty}^{\infty} \prod_{j=1}^{k-1} \Phi \left( \frac{\sigma_k}{\sigma_j} z + \frac{\sqrt{r} \Delta}{\sigma_j} \right) d\Phi(z) \\ &\geq \int_{-\infty}^{\infty} \prod_{j=1}^{k-1} \Phi \left( \delta_j z + \frac{\sqrt{r} \Delta}{\sigma_{[j+1]}} \right) d\Phi(z) \text{ (by a lemma in [4])} \end{aligned}$$

where  $\delta_j = \sigma_{[1]} / \sigma_{[j+1]}$ ,  $\sigma_{[1]} \leq \sigma_{[2]} \leq \dots \leq \sigma_{[k]}$ .

This completes the proof.

Asymptotic relative efficiency of  $R_p$  with respect to a procedure  $R_B$

We assume  $m=2$ ,  $t=1$ , and  $\lambda$  is a projection function. In this case we have  $G_i(x) = \alpha_i F_1(x; \theta_1) + (1 - \alpha_i) F_2(x; \theta_2)$  for  $i=1, 2, \dots, k$  and we denote  $\alpha_i$  instead of  $\alpha_i$ . Suppose  $F_1(x; \theta_1)$  and  $F_2(x; \theta_2)$  are not specified, however, we assume there exists some point  $x_0$ , known, such that  $F_1(x_0; \theta_1) \neq F_2(x_0; \theta_2)$ . Assume  $F_1(x_0; \theta_1) > F_2(x_0; \theta_2)$ . Then, we see that  $\alpha_i > \alpha_j$  if, and only if  $G_i(x_0) > G_j(x_0)$ . Hence, selecting best is equivalent to selecting the population associated with the largest  $G(x_0; \alpha_i)$  value.

For a given  $i$ ,  $1 \leq i \leq k$ , and  $j$ ,  $1 \leq j \leq n$ , define

$$Y_{ij} = \begin{cases} 1 & \text{if } X_{ij} \leq X_0 \\ 0 & \text{otherwise} \end{cases}$$

and define

$$\hat{G}_i(X_0) = \sum_{j=1}^n Y_{ij}.$$

Then, it is obvious that  $\hat{G}_i(X_0)$  is binomial random variable with cdf  $B(n; G(x_0))$ .

We define a selection procedure  $R_B$  as follows:

$R_B$ : Select the population  $\pi_i$  which is associated with the largest  $\hat{G}_i(x_0)$ .

When  $n$  is large, we use the normal approximation. Let  $F_1(X_0; \theta_1) - F_2(X_0; \theta_2) = d_0 > 0$ . Then, by the result of [10], we have, asymptotically  $n \approx c^2(p^*)(1-\Delta^2 d_0^2)/2\Delta^2 d_0^2$ , when  $\Delta \rightarrow 0$ , and  $p^* \rightarrow 1$ . Again, by the Feller's inequality, we see that  $\Phi(z) \approx 1 - \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$ . We obtain thus  $C^2(p^*) = \left(\frac{1}{1-p^*}\right)^2$ . Let  $n_1$  and  $n_2$  denote, respectively, the sample sizes associated with  $R_p$  and  $R_B$  when  $\inf_{\alpha \in \Omega(\lambda; \Delta)} p\{CS\} = p^*$  is satisfied for both rules. We define the asymptotic relative efficiency of  $R_p$  with respect to  $R_B$  by  $ARE(R_p; R_B) = \frac{n_1(p^*, \Delta)}{n_2(p^*, \Delta)}$  as  $p^* \rightarrow 1$  and then  $\Delta \rightarrow 0$ . It follows from the previous result and Theorem 3.2. We have

$$ARE(R_p; R_B) = \lim_{\Delta \rightarrow 0} \lim_{p^* \rightarrow 1} \frac{C\Delta^{\frac{-2}{1-2\delta}}}{1} = 0.$$

$$2d_0^2 \Delta^2 (1-p^*)^2$$

However, if we take  $1-p^* \equiv a = \Delta \rightarrow 0$ , we have our another kind of efficiency given by

$$\text{ARE}(R_p; R_B) = \lim_{\Delta \rightarrow 0} \Delta^{2-\delta} = 0 \quad \text{for } 0 < \delta < 1/2.$$

This shows that  $R_p$  is good compared to  $R_B$ . Also  $R_p$  holds for any general  $m$  and  $t$ . We should note that the case  $m = 2$  and  $m > 2$  are quite different and  $R_B$  is useful only for  $m = 2$ .

### (B) Discrete Case

In this case, we denote  $F_1, F_2, \dots, F_m$  as discrete distribution such that the outcomes from each distribution with cdf  $F_i$ , for some  $i$ , can be classified into  $s$  ( $\geq 2$ ) states. Let the probability that an outcome from  $F_i$  belongs to state  $l$  be denoted by  $p_{il}$ . We assume  $F_1, F_2, \dots, F_m$  are all specified and  $p_{il}$  are all given.

For  $\alpha_i \in \langle 0, 1 \rangle^m$  we define a mixture distribution  $G_i$  by

$$G_i(x) = \alpha_{i1}F_1(x) + \alpha_{i2}F_2(x) + \dots + \alpha_{im}F_m(x).$$

Then,  $G_i(x)$  is also a discrete distribution such that the probability of an outcome belonging to state  $j$  is given by

$$g_{ij} = \alpha_{i1}p_{1j} + \alpha_{i2}p_{2j} + \dots + \alpha_{im}p_{mj}, \quad \text{for } j=1, 2, \dots, s.$$

We assume there exists a lower bound  $g_0$  such that  $g_{ij} \geq g_0 > 0$  for all  $i=1, 2, \dots, k, j=1, 2, \dots, s$ . Let  $n$  samples be drawn from  $\pi_i$  and let  $n_j$  denote the number of outcomes which belong to state  $j$ . For any  $\alpha = (\alpha_1, \dots, \alpha_m)$ , we define the Matusita distance (see [8]) as follows.

$$(3.3) \quad S_i(\alpha) = \left\{ \sum_{j=1}^s \left( \sqrt{g_j} \sqrt{\frac{n_j}{n}} \right)^2 \right\}^{\frac{1}{2}}$$

where  $g_j = \sum_{i=1}^m \alpha_i p_{ij}$ .  $S_i(\alpha)$  is thus a function on  $\langle 0, 1 \rangle^m$ .



Let  $\hat{\alpha}_{\sim i}$  denote a value in  $\langle 0,1 \rangle^m$  such that  $S_i(\hat{\alpha}_{\sim i})$  attains its infimum, i.e. let  $\hat{\alpha}_{\sim i}$  be such that

$$(3.4) \quad S_i(\hat{\alpha}_{\sim i}) = \inf_{\alpha \in \langle 0,1 \rangle^m} S_i(\alpha).$$

For given  $n$  and  $\lambda$ , to select the  $t$  best with respect to  $\lambda$ , we propose the following selection procedure.

R: Select  $\pi_{r_1}, \pi_{r_2}, \dots, \pi_{r_t}$  if, and only if,

$\lambda(\hat{\alpha}_{\sim r_1}), \lambda(\hat{\alpha}_{\sim r_2}), \dots, \lambda(\hat{\alpha}_{\sim r_t})$  are the  $t$  largest values of  $\lambda(\hat{\alpha}_{\sim 1}), \lambda(\hat{\alpha}_{\sim 2}), \dots, \lambda(\hat{\alpha}_{\sim k})$  which are defined by (3.3) and (3.4). If there are ties, use a random mechanism.

Theorem 3.4 The selection procedure R is consistent and strongly asymptotically monotone, with respect to  $(\mathfrak{F}, \lambda)$ .

Proof: It has been shown in Matusita [8] that for our case  $\hat{\alpha}_{\sim i} \rightarrow \alpha_{\sim i}$  with probability one in the usual sense of convergence of a sequence of vectors. Therefore,  $\lambda(\hat{\alpha}_{\sim i}) \rightarrow \lambda(\alpha_{\sim i})$  WP1 for  $\lambda$  is continuous. Using the analogous arguments given in the proofs of Theorem 3.1, we can conclude the same results. This completes the proof.

## REFERENCES

- [1] Bartlett, M.S. and MacDonald, P.D.M. (1968) "Least-squares" estimation of distribution mixtures. Nature, Lond., 217 195-196.
- [2] Choi, W. (1969) Estimates for the parameters of a finite mixture of distributions. Ann. Inst. Statist. Math. 21, 107-116.
- [3] Choi, W. and Bulgren, W. G. (1968) An estimation procedure for mixtures of distributions. J. R. Statist. Soc. Ser. B 30, 444-460.
- [4] Gupta, S. S. and Huang, W. T. (1974) A note on selecting a subset of normal populations with unequal sample sizes. Sankhyā, Ser. A, 36, 389-396.
- [5] Hill, B. M. (1963) Information for estimating the proportions in mixtures of exponential and normal distributions. J. Amer. Statist. Assoc. 58, 918-932.
- [6] Lehmann, E. L. (1951) Consistency and unbiasedness of certain non-parametric tests. Ann. Math. Statist. 22, 165-179.
- [7] Macdonald, P.D.M. (1971) Comment on "An estimation procedure for mixtures of distributions" by Choi and Bulgren. J. R. Statist. Soc. Ser. B, 33, 326-329.
- [8] Matusita, K. (1954) On the estimation by the minimum distance method. Ann. Inst. Statist. Math. 5, 59-65.
- [9] Robbins, H. (1964) The empirical Bayes approach to statistical decision problems. Ann. Math. Statist. 35, 1-20.
- [10] Sobel, M. and Huyett, M. (1957). Selecting the best one of several binomial populations. Bell System. Tech. J. 36, 537-576.
- [11] Teicher, H. (1963) Identifiability of finite mixtures. Ann. Math. Statist., 34, 1265-1269.
- [12] Wolfowitz, J. (1954). Estimation by the minimum distance method. Ann. Inst. of Statist. Math. 9-23.
- [13] Yakowitz, S. and Spragins, J. (1968) On the identifiability of finite mixtures. Ann. Math. Statist., 39, 209-214.

~~Unclassified~~

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

REPORT DOCUMENTATION PAGE		READ INSTRUCTIONS BEFORE COMPLETING FORM
1. REPORT NUMBER Mimeograph Series #504	2. GOVT ACCESSION NO.	3. RECIPIENT'S CATALOG NUMBER
4. TITLE (and Subtitle) On Selection Rules for Finite Mixtures of Distributions		5. TYPE OF REPORT & PERIOD COVERED Technical
		6. PERFORMING ORG. REPORT NUMBER Mimeo. Series #504
7. AUTHOR(s) Shanti S. Gupta and Wen-Tao Huang		8. CONTRACT OR GRANT NUMBER(s) ONR N00014-75C-0455
		9. PERFORMING ORGANIZATION NAME AND ADDRESS Purdue University Department of Statistics West Lafayette, IN 47907
11. CONTROLLING OFFICE NAME AND ADDRESS Office of Naval Research Washington, DC		10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS
		12. REPORT DATE August 1977
14. MONITORING AGENCY NAME & ADDRESS (if different from Controlling Office)		13. NUMBER OF PAGES 17
		15. SECURITY CLASS. (of this report) Unclassified
16. DISTRIBUTION STATEMENT (of this Report)  Approved for public release, distribution unlimited.		15a. DECLASSIFICATION/DOWNGRADING SCHEDULE
		17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)
18. SUPPLEMENTARY NOTES		
19. KEY WORDS (Continue on reverse side if necessary and identify by block number) Finite mixture of distributions, consistent selection procedure, strong asymptotic monotonicity, minimum distance, least squares.		
20. ABSTRACT (Continue on reverse side if necessary and identify by block number) A class of consistent selection procedures is proposed for the selection of the $t$ largest proportions of finite mixture of distributions. The approach is based on the indifference zone formulation. Some asymptotic optimal properties are shown to hold. For special case of $m = 2$ , the number of components in each of $k$ mixture populations, some results of the asymptotic relative efficiency are obtained.		

DD FORM 1473  
1 JAN 73

EDITION OF 1 NOV 65 IS OBSOLETE  
S/N 0103-014-6601

Unclassified

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)