# ON ESTIMATING OF THE NUMBER OF CONSTITUENTS OF A FINITE MIXTURE OF CONTINUOUS DISTRIBUTIONS

Jōgi Henna

## Summary

Suppose that $H$ is a mixture of distributions for a given family $\mathscr{F}$. A necessary and sufficient condition is obtained under which $H$ is, in fact, a finite mixture. An estimator of the number of distributions constituting the mixture is proposed assuming that the mixture is finite and its asymptotic properties are investigated.

## 1. Introduction

Let $\mathscr{F} = \{F_\theta(x): \theta \in R_1^k\}$ be a family of known one-dimensional cumulative distribution functions (cdf's) and $G^\circ(\theta)$ any cdf such that $P_{G^\circ}(R_1^k) = 1$, where $P_{G^\circ}$ is the probability measure induced by $G^\circ$ and $R_1^k$ a compact subset of the $k$-dimensional Euclidean space $R^k$. Let $F_\theta(x)$ be continuous in $x$ for each $\theta$ and continuous in $\theta$ for each $x$. Let $H_{G^\circ}(x)$ be the continuous cdf defined by

$$(1.1) \qquad H_{G^\circ}(x) = \int_{R_1^k} F_\theta(x) dG^\circ(\theta) .$$

$H_{G^\circ}$ will be called a mixture of $\mathscr{F}$ and, especially, a finite mixture when the support of $P_{G^\circ}$ is a finite subset of $R_1^k$.

Suppose that $H_{G^\circ}$ is known to be a finite mixture and let $X = (X_1, X_2, \cdots, X_n)$ be a random sample from $H_{G^\circ}$. Our problem is to construct an estimator $\hat{m}_n$ of the number $m_0$ of cdf's constituting $H_{G^\circ}$ and to investigate its asymptotic property. The method employed here consists of examining the number of cdf's constituting the finite mixture which is "closest" to the empirical distribution function of $X$. Still another important problem discussed in this paper is a criterion as to whether or not a given mixture is finite.

To the author's knowledge ([2], [3] and others), no estimator of $m_0$

---

has been proposed so far. What is then the purpose of estimating $m_0$? For instance, the estimator $\hat{m}_n$ is useful in classification problems. More precisely, we may regard it as the number of populations when we wish to classify an observation into one of several populations.

In Section 2, we give some notation, an estimator $\hat{m}_n$ and a preliminary lemma. In Section 3, we give a necessary and sufficient condition for $H_{G^0}$ to be a finite mixture. In Section 4, we show that $\hat{m}_n = m_0$ holds with probability one for all $n$ sufficiently large. In Section 5, we give some examples.

## 2. Notations and a preliminary lemma

We define a subfamily of mixing cdf's by

$$\mathcal{G}_m = \Big\{ G_m :\ G_m = (g_1, g_2, \cdots, g_m;\ \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_m),$$

$$\sum_{j=1}^{m} g_j = 1,\ 0 \leqq g_j \leqq 1,\ \boldsymbol{\theta}_j \in R_1^k,\ j = 1, 2, \cdots, m \Big\},$$

$$(m = 1, 2 \cdots),$$

where $G_m = (g_1, g_2, \cdots, g_m;\ \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_m)$ is the discrete cdf with jump $g_j$ at $\boldsymbol{\theta}_j$ $(j = 1, 2, \cdots, m)$.

Let $\hat{G}_{m,n}$ be any $G_m$ in $\mathcal{G}_m$ which minimizes

$$S_n(G_m) = \int \{H_{G_m}(x) - F_n(x)\}^2 dF_n(x) = \frac{1}{n} \sum_{i=1}^{n} \Big\{ \sum_{j=1}^{m} g_j F_{\boldsymbol{\theta}_j}(X_{(i)}) - \frac{i}{n} \Big\}^2,$$

where $F_n(x)$ and $X_{(i)}$ are the empirical distribution function and the $i$th order statistic of $X$ respectively. The existence of $\hat{G}_{m,n}$ is guaranteed since $S_n(G_m)$ is a continuous function of $G_m$ on a compact subset of $R^{m(k+1)}$.

For the case $G^0 = G_{m_0}^0$ with finite but unknown $m_0$, where $0 < g_j^0 < 1$ $(j = 1, 2, \cdots, m_0)$, we propose an estimator $\hat{m}_n$ of $m_0$ defined as follows:

$$\hat{m}_n = \text{the minimal integer } m \text{ such that } S_n(\hat{G}_{m,n}) < \lambda^2(n)/n,$$

where $\lambda(n) \uparrow \infty$, $\lambda^2(n)/n \to 0$ as $n \to \infty$ and $\sum \{\lambda^2(n)/n\} e^{-2\lambda^2(n)} < \infty$. The existence of $\hat{m}_n$ for all $n$ sufficiently large is guaranteed with probability one by Lemma 4.3 which will be shown later.

The following lemma is a simple consequence of Polya's theorem (see Rao [4], p. 120) and the Glivenko-Cantelli theorem. We omit the proof.

LEMMA 2.1. *Let* $\{G_{m,n}\}_{n=1}^{\infty}$ *be any sequence of cdf's in* $\mathcal{G}_m$ *such that* $G_{m,n} \to G_m^*$ *as* $n \to \infty$. *Then*

$$\int \{H_{G_{m,n}}(x) - H_{G^\cdot}(x)\}^2 dF_n(x) \to \int \{H_{G_m^*}(x) - H_{G^\cdot}(x)\}^2 dH_{G^\cdot}(x)$$

*with probability one (with respect to $P_{H_{G^\cdot}}^{(\infty)}$).*

For a sequence $\{G_{m,n}\} = \{(g_{1,n}, \cdots, g_{m,n}; \theta_{1,n}, \cdots, \theta_{m,n})\}$, the convergence $G_{m,n} \to G_m^* = (g_1^*, \cdots, g_m^*; \theta_1^*, \cdots, \theta_m^*)$ means $g_{j,n} \to g_j^*$ and $\theta_{j,n} \to \theta_j^*$ ($j=1, 2, \cdots, m$) as $n \to \infty$.

## 3. A necessary and sufficient condition for $H_{G^\cdot}$ to be a finite mixture

In order to give a necessary and sufficient condition for $H_{G^\cdot}$ to be a finite mixture on the basis of $X$, we need the following identifiability condition.

(A-3.1)   For any finite mixture $H_G$, the relationship $H_G = H_{G^\cdot}$ implies that $G = G^*$.

LEMMA 3.1.  *Under (A-3.1), suppose that $H_{G^\cdot}$ is not a finite mixture. Then*

$$P_{H_{G^\cdot}}^{(\infty)}\{\liminf_{n\to\infty} S_n(\hat{G}_{m,n}) > 0\} = 1$$

*for any finite $m$.*

PROOF.  Assume that the conclusion does not hold. Then there exists a Borel subset $A$ of $R^\infty$ such that $P_{H_{G^\cdot}}^{(\infty)}(A) > 0$ and, if $(X_1, X_2, \cdots) \in A$, then $\|H_{G^\cdot} - F_n\| \to 0$ as $n \to \infty$ (by the Glivenko-Cantelli theorem) and $\liminf\limits_{n\to\infty} S_n(\hat{G}_{m,n}) = 0$ for a finite $m$, where $\| \ \|$ denotes the sup norm. Then there exists a subsequence $\{\hat{G}_{m,r}\}$ of $\{\hat{G}_{m,n}\}$ such that

$$(3.1) \qquad S_r(\hat{G}_{m,r}) = \int \{H_{\hat{G}_{m,r}}(x) - F_r(x)\}^2 dF_r(x) \to 0$$

as $r \to \infty$. Let $\{\hat{G}_{m,s}\}$ be any subsequence of $\{\hat{G}_{m,r}\}$ such that $\hat{G}_{m,s} \to G_m^*$ as $s \to \infty$, where $G_m^*$ is a member of $\mathcal{G}_m$. We have

$$(3.2) \qquad \int \{H_{\hat{G}_{m,s}}(x) - F_s(x)\}^2 dF_s(x)$$

$$= \int \{H_{\hat{G}_{m,s}}(x) - H_{G^\cdot}(x)\}^2 dF_s(x)$$

$$+ 2 \int \{H_{\hat{G}_{m,s}}(x) - H_{G^\cdot}(x)\}\{H_{G^\cdot}(x) - F_s(x)\} dF_s(x)$$

$$+ \int \{H_{G^\cdot}(x) - F_s(x)\}^2 dF_s(x) .$$

The second and third terms on the right hand side of (3.2) converges

to 0 as $s \to \infty$ by $\|H_{G^\cdot} - F_s\| \to 0$. Hence, by Lemma 2.1, we have

$$\int \{H_{\hat{G}_{m,s}}(x) - F_s(x)\}^2 dF_s(x) \to \int \{H_{G_m^*}(x) - H_{G^\cdot}(x)\}^2 dH_{G^\cdot}(x) .$$

It follows from (3.1) and the assumption (A-3.1) that $G_m^* = G^\circ$, i.e., that $H_{G^\cdot}$ is a finite mixture, contradicting to the assumption of the lemma.

COROLLARY 3.1. *Under* (A-3.1), *a necessary and sufficient condition for* $H_{G^\cdot}$ *to be a finite mixture is that there exists a finite* $m$ *such that* $S_n(\hat{G}_{m,n}) \to 0$ *with probability one as* $n \to \infty$ (*with respect to* $P_{H_{G^\cdot}}^{(\infty)}$).

PROOF. Assume that $G^\circ = G_{m_\circ}^\circ$ with finite $m_\circ$. Then, by the definition of $\hat{G}_{m_\circ, n}$, we have

$$0 \le S_n(\hat{G}_{m_\circ, n}) \le \int \{H_{G_{m_\circ}^\cdot}(x) - F_n(x)\}^2 dF_n(x) \le \|H_{G_{m_\circ}^\cdot} - F_n\|^2 .$$

Accordingly, $S_n(\hat{G}_{m_\circ, n}) \to 0$ with probability one as $n \to \infty$ by the Glivenko-Cantelli theorem.

If, conversely, $H_{G^\cdot}$ is not a finite mixture, then the probability that $S_n(\hat{G}_{m,n}) \to 0$ as $n \to \infty$ is 0 for any finite $m$ by Lemma 3.1.

## 4. The asymptotic behavior of $\hat{m}_n$ as $n \to \infty$ in the case of a finite mixture

In order to prove that $\hat{m}_n = m_\circ$ holds with probability one for all $n$ sufficiently large in case $G^\circ = G_{m_\circ}^\circ$ with finite $m_\circ$, where $0 < g_j^\circ < 1$ ($j = 1, 2, \cdots, m_\circ$), we assume the following identifiability condition which is weaker than (A-3.1).

(A-4.1)   For any two finite mixture $H_G$ and $H_{G^*}$, the relationship $H_G = H_{G^*}$ implies that $G = G^*$.

To derive the asymptotic behavior of $\hat{m}_n$, we need to investigate some asymptotic properties of $S_n(\hat{G}_{m,n})$ as $n \to \infty$. We can show the following lemma in the same way to the proof of Lemma 3.1.

LEMMA 4.1. *Under* (A-4.1), *suppose that* $m < m_\circ$. *Then*

$$P_{H_{G_{m_\circ}}^\cdot}^{(\infty)} \{\liminf_{n \to \infty} S_n(\hat{G}_{m,n}) > 0\} = 1 .$$

LEMMA 4.2. $S_n(\hat{G}_{m,n}) \ge S_n(\hat{G}_{m+1,n})$ *for any* $n$.

PROOF.

$$S_n(\hat{G}_{m,n}) = \min S_n(g_1, g_2, \cdots, g_m, 0; \theta_1, \theta_2, \cdots, \theta_m, \theta_{m+1})$$

$$\geq \min S_n(g_1, g_2, \cdots, g_m, g_{m+1}; \; \boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \cdots, \boldsymbol{\theta}_m, \boldsymbol{\theta}_{m+1})$$
$$= S_n(\hat{G}_{m+1,n}) \; .$$

LEMMA 4.3. *Suppose that* $m_o \leq m$. *Then*

$$P_{H_{G_{m_\bullet}}}^{(\infty)} \{S_n(\hat{G}_{m,n}) < \lambda^2(n)/n \text{ for all } n \text{ sufficiently large}\} = 1 \; .$$

PROOF. Suppose that $m_o \leq m$. Then, by the last lemma, we have

$$0 \leq S_n(\hat{G}_{m,n}) \leq S_n(\hat{G}_{m_o,n}) \leq \|H_{G_{m_\bullet}} - F_n\|^2 \; .$$

Accordingly, we have the conclusion by Theorem 2 of Chung [1].

THEOREM 4.1. *Under* (A–4.1), *we have*

$$P_{H_{G_{m_\bullet}}}^{(\infty)} \{\hat{m}_n = m_o \text{ for all } n \text{ sufficiently large}\} = 1 \; .$$

PROOF. By Lemmas 4.1 and 4.2, we have

$$P_{H_{G_{m_\bullet}}}^{(\infty)} \{\hat{m}_n \leq m_o - 1 \text{ for an infinite number of } n\text{'s}\}$$
$$\leq P_{H_{G_{m_\bullet}}}^{(\infty)} \{S_n(\hat{G}_{m_o-1,n}) < \lambda^2(n)/n \text{ for an infinite number of } n\text{'s}\}$$
$$= 0 \; .$$

By Lemmas 4.2 and 4.3, we have

$$P_{H_{G_{m_\bullet}}}^{(\infty)} \{\hat{m}_n \geq m_o + 1 \text{ for an infinite number of } n\text{'s}\}$$
$$\leq P_{H_{G_{m_\bullet}}}^{(\infty)} \{S_n(\hat{G}_{m_o,n}) \geq \lambda^2(n)/n \text{ for an infinite number of } n\text{'s}\}$$
$$= 0 \; .$$

## 5. Examples

The following propositions are derived by a modification of the proof of Theorem 2 of Teicher [5].

PROPOSITION 5.1. *Let* $N_{(\theta,\sigma^2)}(x)$ *be the normal distribution function with mean* $\theta$ *and variance* $\sigma^2$. *Suppose that*

(5.1)
$$\int_D N_{(\theta,\sigma^2)}(x) dG_1(\theta, \sigma^2) = \int_D N_{(\theta,\sigma^2)}(x) dG_2(\theta, \sigma^2)$$

*and either side of* (5.1) *is a finite mixture, where* $P_{G_1}(D) = P_{G_2}(D) = 1$ *and* $D = (\theta_o, \infty) \times (0, \infty)$ *with finite* $\theta_o$. *Then* $G_1 = G_2$.

PROPOSITION 5.2. *Let* $F_{(\alpha,\beta)}(x) = \Gamma(\alpha)^{-1} \int_\beta^x (y - \beta)^{\alpha-1} e^{-(y-\beta)} dy$. *Suppose that*

(5.2)
$$\int_D F_{(\alpha,\beta)}(x) dG_1(\alpha, \beta) = \int_D F_{(\alpha,\beta)}(x) dG_2(\alpha, \beta)$$

and either side of (5.2) is a finite mixture, where $P_{G_1}(D)=P_{G_2}(D)=1$ and $D=(0, \infty) \times (-\infty, \infty)$. Then $G_1=G_2$.

Now we shall give some examples.

*Example* 5.1. Let $\mathscr{F}=\{N_{(\theta, \sigma^2)}(x): (\theta, \sigma^2) \in R_1^2\}$, where $R_1^2$ is a compact subset of $(-\infty, \infty) \times (0, \infty)$. The condition (A-3.1) (accordingly (A-4.1)) is satisfied by Proposition 5.1. So, Corollary 3.1 can be applied. Let

$$\hat{m}_n = \text{the minimal integer } m \text{ such that } S_n(\hat{G}_{m,n}) < \frac{(\log n)^2}{n}.$$

Then Theorem 4.1 can be applied, that is,

$$P_{H_{G^*_{m_o}}}^{(\infty)}\{\hat{m}_n = m_o \text{ for all } n \text{ sufficiently large}\} = 1.$$

*Example* 5.2. Let $\mathscr{F}=\{F_{(\alpha, \beta)}(x); (\alpha, \beta) \in R_1^2\}$, where $R_1^2$ is a compact subset of $(0, \infty) \times (-\infty, \infty)$. The condition (A-3.1) (accordingly (A-4.1)) is satisfied by Proposition 5.2. So, Corollary 3.1 can be applied. Let $\hat{m}_n$ be that of the last example. Then Theorem 4.1 can be applied.

*Example* 5.3. Let $\mathscr{F}=\{F_{(\theta, \alpha)}(x): (\theta, \alpha) \in R_1^2\}$, where $F_{(\theta, \alpha)}(x)=\theta^\alpha[\Gamma(\alpha)]^{-1}$ $\cdot \int_0^x y^{\alpha-1}e^{-\theta y}dy$ and $R_1^2$ is a compact subset of $(0, \infty) \times (0, \infty)$. The condition (A-4.1) is satisfied by Proposition 2 of Teicher [5]. Let $\hat{m}_n$ be that of the last example. Then Theorem 4.1 can be applied.

## Acknowledgements

RYUKYU UNIVERSITY

## REFERENCES

[1] Chung, K. L. (1949). An estimate concerning the Kolmogoroff limit distribution, *Tran. Amer. Math. Soc.*, **67**, 36–50.

[2] Gupta, S. S. and Huang, W. T. (1981). On mixtures of distributions: a survey and some new results on ranking and selection, *Sankhyā*, **43**, 245–290.

[3] Isaenko, O. K. and Urbakh, V. Y. (1977). Partitioning mixed probability distributions into their constituents, *J. Soviet Math.*, **7**, 148–160.

[4] Rao, C. R. (1973). *Linear Statistical Inference and Its Applications*, John Wiley.

[5] Teicher, H. (1963). Identifiability of finite mixtures, *Ann. Math. Statist.*, **34**, 1265–1269.