

AN ABSTRACT OF THE THESIS OF

CHARLES JAMES PHILLIPS for the DOCTOR OF PHILOSOPHY  
(Name) (Degree)

in STATISTICS presented on April 9, 1969  
(Major) (Date)

Title: CONSISTENT EMPIRICAL APPROXIMATION OF A-PRIORI  
DISTRIBUTIONS

Abstract approved: *Redacted for Privacy*  
Donald Guthrie, Jr.

Suppose we have a repeated decision problem based on the random variable  $X$  whose distribution has density  $f(x|\lambda) \in \mathcal{F}$ , where the parameter  $\lambda$  has an unknown cumulative distribution function  $G(\lambda) \in \mathcal{G}$ . The unconditional density of  $X$  is a  $G$ -mixture over  $\mathcal{F}$ ,  $f_G(x) = \int_{\Lambda} f(x|\lambda) dG(\lambda)$ . The estimation of  $G(\lambda)$  using successive observations of the random variable  $X$  is developed in this thesis.

To get approximate estimates of  $G$  for the identifiable subfamily of the exponential family of distributions, a certain function of  $\lambda$  is approximated by a set of functions derived from the parameterization of the exponential family and a sequence of consistent estimators of the expectation of this function is developed.

In the case where  $G$  is a discrete distribution over a finite number of points, we derive consistent estimators based on a sequence of samples whose sizes depend on previously observed

samples. Application to acceptance sampling is indicated. In this latter case the family  $\mathcal{F}$  is enlarged to include all densities which are square integrable and bounded almost everywhere with respect to a suitable measure.

Consistent Empirical Approximation of  
A-priori Distributions

by

Charles James Phillips

A THESIS

submitted to

Oregon State University

in partial fulfillment of  
the requirements for the  
degree of

Doctor of Philosophy

June 1969

APPROVED:

*Redacted for Privacy*

---

Associate Professor of Mathematics and Statistics

in charge of major

*Redacted for Privacy*

---

Chairman of Department of Statistics

*Redacted for Privacy*

---

Dean of Graduate School

Date thesis is presented April 9, 1969

Typed by Clover Redfern for Charles James Phillips

## ACKNOWLEDGMENT

The author wishes to express his thanks to his wife, Jeanette, for her patience and forbearance during the long travail, to his major professor, Dr. Donald Guthrie, Jr., without whose confidence and generous help this thesis would not have been possible, and to the National Science Foundation for support of this research under grants GP-6142 and GP-7491.

## TABLE OF CONTENTS

Chapter	Page
I. INTRODUCTION	1
II. ESTIMATION OF A-PRIORI DISTRIBUTIONS OVER COMPACT SUBSETS OF THE PARAMETER SPACE $\Lambda$	11
III. ESTIMATION OF DISCRETE PRIOR DISTRIBUTIONS WITH VARYING DIMENSIONAL RANDOM VECTORS	36
BIBLIOGRAPHY	48

# CONSISTENT EMPIRICAL APPROXIMATION OF A-PRIORI DISTRIBUTIONS

## I. INTRODUCTION

In the traditional formulation of statistical decision problems the probability distribution of the observed random variable is assumed to be a member of some class  $\{F(x|\lambda), \lambda \in \Lambda\}$  of distribution functions. No prior information is usually assumed to be known concerning which element of the index set is the true one in any given situation. It is possible, of course, to introduce a priori probability measures defined on  $\Lambda$  as a technical device for generating complete classes of decision functions, minimax decision rules, etc., and in some cases as subjective measures of prior information. It would seem quite reasonable, however, that we may assume in certain experimental situations that such an a priori probability measure actually exists on  $\Lambda$  in the sense that the distributions of the observed random variables occurring in different experiments are selected according to some probability distribution  $G(\lambda)$  defined on the index set  $\Lambda$ .

To make such an assumption seems particularly proper in the case where one wishes to make a statistical decision about some characteristic or characteristics of a population at each step of a sequence of independent observations of populations, the decision at

each step depending on some population parameter of the particular population occurring at that step. If the probability distribution of the parameter is known to the investigator then an optimum Bayes decision procedure can be determined. Usually, though, such information will not be available to the investigator, but previous observations on the populations may be and under the proper circumstances these observations may be used to approximate the optimum Bayes procedure. The possibility of using previous observations in such approximations was first conjectured by H. Robbins (1956).

In this paper then we will discuss decision problems arising with the following structure:

- a) An observable random variable  $x$  from an  $\sigma$ -finite measure space  $(X, \chi, \mu)$ .
- b) The random variable  $x$  has a probability density function  $f_{\lambda}(x)$  depending on a parameter  $\lambda$ .
- c) The parameter  $\lambda$  is an element of a space  $\Lambda$  on which is defined a known or unknown a-priori distribution  $G(\lambda)$ .

The specific value of  $\lambda$  in any situation is unknown to us.

- d) An action space  $A$  with generic element  $a$ .
- e) A loss function  $L(a, \lambda) \geq 0$  which represents the loss incurred in taking action  $a$  when the parameter value is  $\lambda$ .
- f) A decision function space  $T$  such that for  $t \in T$ ;  $t : X \rightarrow A$ .

The problem will be to determine a decision function  $t \in T$



such that upon observing  $x$  we take action  $t(x)$  and incur loss  $L(t(x), \lambda)$ . Thus for any  $t$  the expected loss for a given  $\lambda$  is

$$(1.1) \quad R(t, \lambda) = \int_{\mathbf{X}} L(t(\mathbf{x}), \lambda) f_{\lambda}(\mathbf{x}) d\mu(\mathbf{x})$$

which leads to the Bayes risk

$$(1.2) \quad R(t, G) = \int_{\Lambda} R(t, \lambda) dG(\lambda).$$

If we set

$$(1.3) \quad \phi_G(a, \mathbf{x}) = \int_{\Lambda} L(a, \lambda) f_{\lambda}(\mathbf{x}) dG(\lambda)$$

we may then write the Bayes risk as

$$(1.4) \quad R(t, G) = \int_{\mathbf{X}} \phi_G(t(\mathbf{x}), \mathbf{x}) d\mu(\mathbf{x}).$$

For simplicity we assume the existence of the Bayes decision function  $t_G$  such that a. e.  $\mu(\mathbf{x})$

$$(1.5) \quad \phi_G(t_G(\mathbf{x}), \mathbf{x}) = \min_{\mathbf{A}} \phi_G(a, \mathbf{x}).$$

Then for any  $t$

$$(1.6) \quad R(t_G, G) = \int_X \min_A \phi_G(a, x) d\mu(x) \leq R(t, G)$$

and by defining

$$(1.7) \quad R(G) = R(t_G, G) = \int_X \phi_G(t_G(x), x) d\mu(x)$$

we obtain

$$(1.8) \quad R(G) = \min_T R(t, G).$$

As was pointed out above under the proper circumstances we may assume the existence of a distribution  $G(\lambda)$  on  $\Lambda$  where we are faced repeatedly and independently with the above type decision problem, the value of  $\lambda$  occurring according to  $G(\lambda)$ . This problem was first considered by Robbins (1956) in a more restricted form and in the more general form again by Robbins (1963). The non-parametric case was considered by Johns (1956) and the parametric case in a paper slightly less general than Robbins' by Samuel (1963).

In this more general case we have a sequence of pairs of random variables  $(\lambda_1, x_1), (\lambda_2, x_2), \dots$ , each pair independent of all others, the  $\lambda_n$  having the common a-priori distribution  $G$ , and the conditional density of  $x_n$  given  $\lambda_n$  being  $f_{\lambda_n}(x) = f(x|\lambda = \lambda_n)$ . Given the observation  $x_1, x_2, \dots, x_n$  (the values  $\lambda_i$  remaining unknown) we wish to make a decision about the value  $\lambda_{n+1}$  when

observation  $x_{n+1}$  is obtained. Thus we use the decision function

$$(1.9) \quad t_n(\cdot) = t_n(x_1, x_2, \dots, x_n; \cdot),$$

i. e. , we take action  $t_n(x_{n+1})$  and incur the loss  $L(t_n(x_{n+1}), \lambda_{n+1})$ .

Due to this Robbins has defined an "empirical decision procedure" to be a sequence  $T^* = \{t_n\}$  of functions of the form (1.9) with values in  $A$ . We then have the expected loss

$$(1.10) \quad \int_X \phi_G(t_n(x), x) d\mu(x)$$

and over all expected loss

$$(1.11) \quad R_n(T^*, G) = \int_X E \phi_G(t_n(x), x) d\mu(x)$$

with expectation taken w. r. t.  $x_1, \dots, x_n$  where

$$(1.12) \quad f_G(x) = \int_{\Lambda} f_{\lambda}(x) dG(\lambda)$$

and clearly

$$(1.13) \quad R_n(T^*, G) \geq R(G).$$

Under this formulation then there is a sequence of risk functions  $\{R_n\}$  so the following definition is made: If  $\lim_{n \rightarrow \infty} R_n(T^*, G) = R(G)$

then we say  $T^*$  is asymptotically optimal relative to  $G$ . Let

$$(1.14) \quad \Delta_G(a, x) = \int_{\Lambda} [L(a, \lambda) - L(a_0, \lambda)] f_{\lambda}(x) dG(\lambda)$$

for a fixed arbitrary  $a_0 \in A$  and

$$(1.15) \quad L_0(x) = \int_{\Lambda} L(a_0, \lambda) f_{\lambda}(x) dG(\lambda).$$

Then if

$$\int_{\Lambda} L(\lambda) dG(\lambda) < \infty, \quad L(\lambda) = \sup_A L(a, \lambda)$$

$$(1.16) \quad \phi_G(a, x) = L_0(x) + \Delta_G(a, x), \quad \text{a. e. } \mu(x).$$

If we have a sequence of functions

$$(1.17) \quad \Delta_n(a, x) = \Delta_n(x_1, \dots, x_n; a, x)$$

such that a. e.  $\mu(x)$

$$(1.18) \quad p \lim_{n \rightarrow \infty} \sup_A |\Delta_n(a, x) - \Delta_G(a, x)| = 0$$

and a sequence

$$(1.19) \quad t_n(x) = t_n(x_1, \dots, x_n; x) = \text{any element } \bar{a} \text{ of } A \text{ such that}$$

$$\Delta_n(\bar{a}, x) \leq \inf_A \Delta_n(a, x) + \epsilon_n \quad \text{where } \{\epsilon_n\} \rightarrow 0$$

then by Theorem 1 of Robbins paper (1963)  $T^* = \{t_n\}$  is said to be asymptotically optimal relative to  $G$ .

If we restrict  $A$  to the case  $A = \{a_0, a_1\}$ , we may proceed as follows. Find a sequence  $G_n(\lambda) = G_n(x_1, x_2, \dots, x_n; \lambda)$  of random distribution functions in  $\lambda$  such that

$$(1.20) \quad P \left[ \lim_{n \rightarrow \infty} G_n(\lambda) = G(\lambda) \text{ at every continuity point of } G \right] = 1.$$

One immediately notices that we are then proceeding with a sequence of estimators  $G_n(\lambda)$  which is a consistent estimator of  $G(\lambda)$ , i. e., pointwise consistency at continuity points of  $G(\lambda)$ . Set

$$(1.21) \quad \Delta_n(x) = \int_{-\infty}^{\infty} [L(a_1, \lambda) - L(a_0, \lambda)] f_{\lambda}(x) dG_n(\lambda)$$

and if a. e.  $\mu(x)$ , fixed  $x$ ,

$$(1.22) \quad [L(a_1, \lambda) - L(a_0, \lambda)] f_{\lambda}(x)$$

is continuous and bounded in  $\lambda$  then

$$(1.23) \quad p \lim_{n \rightarrow \infty} \Delta_n(x) = \Delta_G(x) = \int_{-\infty}^{\infty} [L(a_1, \lambda) - L(a_0, \lambda)] f_{\lambda}(x) dG(\lambda)$$

which by (1.19) gives us asymptotic optimality of  $T^*$  relative to  $G$ .

Thus we see that a pointwise consistent (at points of continuity of  $G(\lambda)$ ) sequence of pointwise estimators of  $G(\lambda)$  produces asymptotic optimality of  $T^*$  relative to  $G$ .

In the construction of a sequence of consistent estimators we will need the concept of identifiability of mixtures. There are two particularly good papers on this topic by Teicher (1960, 1963). Let  $\mathcal{F} = \{f(x|\lambda), \lambda \in \Lambda\}$  where  $f(x|\lambda)$  is a density function in the random variable  $x$  for each  $\lambda \in \Lambda$  and measurable on the product space  $X \times \Lambda$ . Then for any non-degenerate (c.d.f.  $G$ ) such that

$$(1.24) \quad f_G(x) = \int_{\Lambda} f(x|\lambda) dG(\lambda)$$

is a density function of the random variable  $x$ ,  $f_G(x)$  is called a  $G$ -mixture of  $\mathcal{F}$  or, simply, a mixture.

Let  $\mathcal{L}$  denote a class of such c.d. f's  $\{G\}$ ,  $\mathcal{H}$  the induced class of mixtures  $\{f\}$  and  $\mathcal{D}$  the class of degenerate distributions in  $\Lambda$ . Then  $\mathcal{H}$  will be called identifiable in  $\mathcal{L}$  (with respect to  $\mathcal{F}$ ) if (1.24) effects a one-to-one correspondence between  $\mathcal{H} \cup \mathcal{F}$  and  $\mathcal{L} \cup \mathcal{D}$ ; equivalently, if

$$f(x) = \int_{\Lambda} f(x|\lambda) dG_1(\lambda) = \int_{\Lambda} f(x|\lambda) dG_2(\lambda)$$

implies  $G_1 = G_2$  for all  $G_1, G_2 \in \mathcal{G} \cup \mathcal{D}$ . If  $\mathcal{H}$  is identifiable in the class of all  $G \in \mathcal{D}$ , it is simply called identifiable.

In particular, Teicher shows that the class of mixtures generated by the normal distribution with either the mean or variance fixed is identifiable and the class of gamma mixtures with one of the parameters fixed also identifiable.

Another pair of concepts to be used are those of completeness and closure of sequences in a normed linear space. Given a normed linear space  $X$  we say a sequence  $\{x_n\}$  of elements of  $X$  is closed if every element of  $X$  can be approximated arbitrarily closely by finite linear combinations of the elements of the sequence, where the degree of approximation is measured in terms of the norm of the space. We also say that a sequence  $\{x_n\}$  is complete if  $L(x_n) = 0, n = 0, 1, \dots, L \in X^*$ , implies  $L = 0$ . Here of course,  $L$  is a linear functional in the conjugate space of  $X$ . The fundamental relationship here is the equivalence of these two concepts in a normed linear space. For an excellent discussion of these concepts and their relationships and applications to approximation theory see Davis (1963).

We will address ourselves to two problems in this paper, the first being the construction of a sequence of consistent estimates of  $G(\lambda)$ , considered in Chapter II. Some work has been done on this for finite mixtures by Robbins (1964), and by Rolph (1968) for

mixtures with discrete domain. We will initially assume the family  $\{f_G\}$  of  $G$  mixtures.

$$(1.25) \quad f_G(\mathbf{x}) = \int_{\Lambda} f(\mathbf{x}|\lambda) dG(\lambda)$$

is identifiable. We then construct a sequence of pointwise estimators  $\{G_n(\lambda)\}$  of  $G(\lambda)$ , where  $f(\mathbf{x}|\lambda)$  is a member of the exponential family, by a method of approximation in a normed linear space of certain functions of the parameter. We use the strong law of large numbers to establish consistency.

In Chapter III we discuss a generalization of Robbins's (1964) work on discrete distributions to the case of variable dimensional random vectors with applications to the theory of acceptance sampling.



## II. ESTIMATION OF A-PRIORI DISTRIBUTIONS OVER COMPACT SUBSETS OF THE PARAMETER SPACE $\Lambda$

Suppose that we have repeated observations of the random variable distributed with density

$$(2.1) \quad f_G(\mathbf{x}) = \int_{\Lambda} f(\mathbf{x}|\lambda) dG(\lambda),$$

where

$$(2.2) \quad f(\mathbf{x}|\lambda) \in \mathcal{F}$$

and the parameter  $\lambda$  has the unknown cumulative distribution function

$$(2.3) \quad G(\lambda) \in \mathcal{G}.$$

The problem is, of course, to estimate  $G(\lambda)$  using observations of the random variable  $X$ . We will assume a measure space  $(X, \chi, \mu)$  and cumulative distribution function

$$(2.4) \quad F(\mathbf{x}|\lambda) = P_{\lambda}(X \leq \mathbf{x})$$

which is absolutely continuous with respect to  $\mu$  and

$$(2.5) \quad \int_X f^2(\mathbf{x}|\lambda) d\mu(\mathbf{x}) < \infty, \quad \lambda \in \Lambda.$$

We further assume that the family  $\mathcal{F}$  is identifiable in  $\mathcal{G}$ , i.e.,

$$(2.6) \quad \int_{\Lambda} f(x|\lambda) dG_1(\lambda) = \int_{\Lambda} f(x|\lambda) dG_2(\lambda)$$

implies

$$G_1(\lambda) = G_2(\lambda), \quad f \in \mathcal{F}, \quad G_i \in \mathcal{G}.$$

Suppose further we assume there exists a compact subset  $L$  of  $\Lambda$  such that  $P(\lambda \in \tilde{L}) = 0$ . For our explicit purpose  $\Lambda$  will be the real line and  $L$  a closed interval  $[\lambda_\ell, \lambda_r]$ ,  $-\infty < \lambda_\ell < \lambda_r < \infty$ .

We estimate  $G(\lambda)$  and other expressions involving  $\lambda$  using linear combinations of  $f(x|\lambda_i)$ . Let

$$(2.7) \quad \Phi_n(x|\lambda^*) = \sum_{i=1}^{k_n} a_{i\lambda^*}^{(n)} f(x|\lambda_i^{(n)})$$

where

$$\lambda_i^{(n)} \in \Lambda, \quad i = 1, \dots, k_n, \quad n \in \mathbb{I}^+,$$

and

$$(2.8) \quad \bar{G}_m^n(\lambda^*) = \frac{1}{m} \sum_{\nu=1}^m \Phi_n(x_\nu|\lambda^*).$$

Consider

$$(2.9) \quad E_{\mathbf{X}} \bar{G}_m^n(x|\lambda^*) = \int_{\mathbf{X}} \Phi_n(x|\lambda^*) \int_{\Lambda} f(x|\lambda) dG(\lambda) d\mu(x).$$

In (2.1) we assumed

$$\begin{aligned}
 (2.10) \quad 1 &= \int_{\mathbf{X}} f_G(\mathbf{x}) d\mu(\mathbf{x}) = \int_{\mathbf{X}} \int_{\Lambda} f(\mathbf{x}|\lambda) dG(\lambda) d\mu(\mathbf{x}) \\
 &= \int_{\mathbf{Y}} f(\mathbf{x}|\lambda) d\pi(\mathbf{Y})
 \end{aligned}$$

where  $(\mathbf{Y}, \mathbf{y}, \pi)$  is the measure space generated by  $\mathbf{X} \times \Lambda$ .

Clearly  $f(\mathbf{x}|\lambda_i^{(n)})f(\mathbf{x}|\lambda)$  is  $\pi$  integrable. Thus by Fubini's theorem

$$\begin{aligned}
 (2.11) \quad E_{\mathbf{X}} \Phi_n(\mathbf{x}|\lambda^*) &= \int_{\Lambda} \int_{\mathbf{X}} \Phi_n(\mathbf{x}|\lambda^*) f(\mathbf{x}|\lambda) d\mu(\mathbf{x}) dG(\lambda) \\
 &= \int_{\Lambda} H_n(\lambda|\lambda^*) dG(\lambda)
 \end{aligned}$$

where we let

$$(2.12) \quad \int_{\mathbf{X}} \Phi_n(\mathbf{X}|\lambda^*) f(\mathbf{x}|\lambda) d\mu(\mathbf{x}) = H_n(\lambda|\lambda^*).$$

If

$$H_n(\lambda|\lambda^*) = \begin{cases} 1, & \lambda \leq \lambda^* \\ 0, & \lambda > \lambda^* \end{cases}$$

then

$$E_{\mathbf{X}} \Phi_n(\mathbf{x}|\lambda^*) = G(\lambda^*).$$

Thus, we wish to choose the  $a_{i\lambda^*}^{(n)}$  such that

$$(2.13) \quad H_n(\lambda | \lambda^*) \doteq H(\lambda | \lambda^*)$$

where

$$(2.14) \quad H(\lambda | \lambda^*) = \begin{cases} 1, & \lambda \leq \lambda^* \\ 0, & \lambda > \lambda^* \end{cases} .$$

(Clearly, in general we would not be able to choose the  $a_{i\lambda^*}^{(n)}$  so that  $H_n(\lambda | \lambda^*) = H(\lambda | \lambda^*)$  for all values of  $\lambda$ .) Now

$$(2.15) \quad \begin{aligned} H_n(\lambda | \lambda^*) &= \int_{\mathbf{X}} \sum_{i=1}^{k_n} a_{i\lambda^*}^{(n)} f(\mathbf{x} | \lambda_i^{(n)}) f(\mathbf{x} | \lambda) d\mu(\mathbf{x}) \\ &= \sum_{i=1}^{k_n} a_{i\lambda^*}^{(n)} \int_{\mathbf{X}} f(\mathbf{x} | \lambda_i^{(n)}) f(\mathbf{x} | \lambda) d\mu(\mathbf{x}) . \end{aligned}$$

Let

$$(2.16) \quad \int_{\mathbf{X}} f(\mathbf{x} | \lambda_i^{(n)}) f(\mathbf{x} | \lambda) d\mu(\mathbf{x}) = h_{\lambda_i^{(n)}}(\lambda)$$

so that

$$(2.17) \quad H_n(\lambda | \lambda^*) = \sum_{i=1}^{k_n} a_{i\lambda^*}^{(n)} h_{\lambda_i^{(n)}}(\lambda) .$$

We have the following

Lemma 1.  $\{f(\mathbf{x} | \lambda_i^{(n)}), \dots, f(\mathbf{x} | \lambda_{k_n}^{(n)})\}$

is a linearly independent set of functions for  $\lambda_i^{(n)} \neq \lambda_j^{(n)}$ ,  $i \neq j$ .

Proof. Suppose  $C_1 f_1^{(n)} + C_2 f_2^{(n)} + \dots + C_{k_n} f_{k_n}^{(n)} = 0$  for

$$\sum_{i=1}^{k_n} C_i^2 \neq 0.$$

Then we may rearrange and renumber so that

$$C_1 f_1^{(n)} + \dots + C_\ell f_\ell^{(n)} = C_{\ell+1} f_{\ell+1}^{(n)} + \dots + C_q f_q^{(n)},$$

where  $C_i > 0 \forall_i \leq q$ . Then integrating with respect to  $\mu(x)$  we have

$$C_1 + C_2 + \dots + C_\ell = C_{\ell+1} + \dots + C_q = C > 0$$

and

$$\frac{C_1}{C} + \dots + \frac{C_\ell}{C} = \frac{C_{\ell+1}}{C} + \dots + \frac{C_q}{C} = 1.$$

Let  $G_1(\lambda)$  assign probability  $\frac{C_i}{C}$  to  $f(x|\lambda)$  if  $\lambda = \lambda_i^{(n)}$ ,  $i = 1, 2, \dots, \ell$  and 0 to  $f(x|\lambda)$  otherwise;  $G_2(\lambda)$  assigns probability  $\frac{C_i}{C}$  to  $f(x|\lambda)$  if  $\lambda = \lambda_i^{(n)}$ ,  $i = \ell+1, \dots, q$  and 0 otherwise. Then

$$f_{G_1}(x) = \sum_{i=1}^{\ell} f(x|\lambda_i^{(n)}) \frac{C_i}{C} = \sum_{i=\ell+1}^q f(x|\lambda_i^{(n)}) \frac{C_i}{C} = f_{G_2}(x)$$

which by (2.6) implies  $G_1(\lambda) = G_2(\lambda)$ , and thus we have a contradiction. Then  $C_i = 0$ ,  $\forall_i$  and the lemma is proved.

We immediately have

Lemma 2. The set of functions  $\{h_{\lambda_1}^{(n)}(\lambda), \dots, h_{\lambda_{k_n}}^{(n)}(\lambda)\}$  is a linearly independent set of functions.

Proof. Suppose not. Then

$$\sum_{i=1}^{k_n} b_i^{(n)} h_{\lambda_i}^{(n)}(\lambda) = 0,$$

where

$$\sum_{i=1}^{k_n} (b_i^{(n)})^2 \neq 0$$

thus

$$\int_X \sum_{i=1}^{k_n} b_i^{(n)} f(x|\lambda_i^{(n)}) f(x|\lambda) d\mu(x) = 0$$

This says

$$\sum_{i=1}^{k_n} b_i^{(n)} f(x|\lambda_i^{(n)}) = 0$$

which contradicts the preceding lemma. Thus the lemma is proved.

Let us now restrict ourselves to the subset of the exponential

family which is identifiable. We shall use the representation where the sigma-finite measure  $\mu(x)$  determines the density function of a member of the family in the form

$$(2.18) \quad f(x|\lambda) = \beta(\lambda)e^{\lambda x}$$

where

$$(2.19) \quad \beta(\lambda) = \left( \int_X e^{\lambda x} d\mu(x) \right)^{-1}$$

Lehmann (1959) has shown that  $\beta(\lambda)$  is a continuous and differentiable function in  $\lambda$ , and thus that

$$(2.20) \quad h_{\lambda_i^{(n)}}^{(n)}(\lambda) = \frac{\beta(\lambda_i^{(n)})\beta(\lambda)}{\beta(\lambda_i^{(n)} + \lambda)}$$

is a continuous function on any finite closed interval  $[\lambda_i^{(n)}, \lambda_{k_n}^{(n)}]$ .

Thus

$$(2.21) \quad \int_{\lambda_1^{(n)}}^{\lambda_{k_n}^{(n)}} [h_{\lambda_i^{(n)}}^{(n)}(\lambda)]^2 d\lambda < \infty, \quad (\text{Lebesgue measure}),$$

and the set of  $h_{\lambda_i^{(n)}}^{(n)}(\lambda)$  generates a finite dimensional subspace of the Hilbert space  $H$  of square integrable functions on  $[\lambda_1^{(n)}, \lambda_{k_n}^{(n)}]$ .

Clearly  $H(\lambda|\lambda^*) \in H$ .

Let us consider the case of

$$(2.22) \quad f(x|\lambda) = \frac{1}{\sqrt{2\pi}} e^{-\lambda^2 x^2 / 2}, \quad -\infty < \lambda < \infty,$$

where

$$(2.23) \quad \frac{d\mu}{dm}(x) = e^{-x^2/2},$$

where  $x$  is real and  $m$  denotes Lebesgue measure.

It is shown by Waterman (1966) that this represents the normal distribution with mean  $\lambda$  and variance 1. Immediately we observe that

$$(2.24) \quad \beta(\lambda) = \frac{1}{\sqrt{2\pi}} e^{-\lambda^2/2}$$

so that

$$(2.25) \quad \frac{\beta(\lambda_i^{(n)})\beta(\lambda)}{\beta(\lambda + \lambda_i^{(n)})} = \frac{1}{\sqrt{2\pi}} e^{-\lambda_i^{(n)}\lambda}.$$

If one wishes to consider the case where the variance is  $\sigma^2 \neq 1$ , but fixed, then merely define  $\mu(x)$  by

$$\frac{d\mu}{dm}(x) = e^{-x^2/2\sigma^2},$$

so that

$$\beta(\lambda) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\lambda^2\sigma^2/2}$$



and

$$\frac{\beta(\lambda_i^{(n)})\beta(\lambda)}{\beta(\lambda+\lambda_i^{(n)})} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{\lambda_i^{(n)}}{\sigma^2} \lambda}.$$

This change does not effect the subsequent analysis so it will not be pursued further. Thus let us look at

$$(2.26) \quad \frac{1}{\sqrt{2\pi}} e^{-\frac{\lambda_i^{(n)}}{\sigma^2} \lambda},$$

where

$$(2.27) \quad \lambda_\ell < \lambda_r, \quad \lambda_\ell \leq \lambda \leq \lambda_r, \quad \lambda_\ell \leq \lambda_i^{(n)} \leq \lambda_r,$$

and

$$|\lambda_i^{(n)}| > 0.$$

Let

$$(2.28) \quad \lambda = (\lambda_r - \lambda_\ell)\omega + \lambda_\ell.$$

Then

$$0 \leq \omega \leq 1 \quad \text{for} \quad \lambda_\ell \leq \lambda \leq \lambda_r$$

and

$$(2.29) \quad \begin{aligned} \frac{1}{\sqrt{2\pi}} e^{-\frac{\lambda_i^{(n)}}{\sigma^2} \lambda} &= \frac{1}{\sqrt{2\pi}} e^{-\frac{\lambda_i^{(n)}}{\sigma^2} (\lambda_r - \lambda_\ell)\omega + \frac{\lambda_i^{(n)}}{\sigma^2} \lambda_\ell} \\ &= \frac{1}{\sqrt{2\pi}} e^{-\frac{\lambda_i^{(n)}}{\sigma^2} \lambda_\ell} e^{-\frac{\lambda_i^{(n)}}{\sigma^2} (\lambda_r - \lambda_\ell)\omega}. \end{aligned}$$

In any linear combination of the above terms the constants

$e^{\lambda_i^{(n)} z}$  may be absorbed into the coefficients of  $e^{\lambda_i^{(n)}(\lambda_r - \lambda_\ell) z}$  so that without loss of generality we consider the independent set  $\{e^{\lambda_i^{(n)*} z}\}$

where

$$(2.30) \quad \lambda_i^{(n)*} = \lambda_i^{(n)}(\lambda_r - \lambda_\ell).$$

The following theorem by Szász is proven in Davis (1963):

Let

$$F(z) = \sum_{k=0}^{\infty} C_k z^k, \quad C_k \text{ real,}$$

be a fixed power series and let  $r > 0$  be its radius of convergence. Assume that

$$\sum_{n=1}^{\infty} \frac{1}{k_n} = \infty$$

where  $k_1 < k_2 < \dots$  is the sequence of all those integers  $\geq 1$  for which  $C_k \neq 0$ . If  $\{t_n\}$  is a sequence of distinct real numbers satisfying  $0 < |t_n| \leq r_1 < r$ , then the sequence of functions  $\{f_n(x)\} = \{F(t_n x)\}$ ,  $n = 1, 2, \dots$ , is complete in  $L^2[0, 1]$ . If  $C_0 \neq 0$ , it is also complete in  $C[0, 1]$ .

Now

$$e^z = 1 + \sum_{k=1}^{\infty} \frac{1}{k!} z^k$$

with

$$r = \infty, \quad \text{i. e.}, \quad |z| < \infty, \quad C_0 = 1, \quad \text{and} \quad C_k = \frac{1}{k!}.$$

Thus

$$k_n = n \quad \text{and} \quad \sum_{n=1}^{\infty} \frac{1}{k_n} = \sum_{n=1}^{\infty} \frac{1}{n} = \infty.$$

Now

$$(2.31) \quad 0 < |\lambda_\ell(\lambda_r - \lambda_\ell)| \leq |\lambda_i^{(n)*}| < |\lambda_r(\lambda_r - \lambda_\ell)| < \infty.$$

So let  $i = n$ ,  $k_n = n$ ,

$$(2.32) \quad \lambda_n^{(n)*} = \lambda_\ell(\lambda_r - \lambda_\ell) + \frac{n-1}{n}(\lambda_r - \lambda_\ell)^2, \quad n = 1, 2, \dots$$

Then by Szàsz' theorem  $\{e^{\lambda_n^{(n)*} \omega}\}$ ,  $n = 1, 2, \dots$ , is complete in  $L^2[0, 1]$  and also in  $C[0, 1]$ . Thus with

$$(2.33) \quad \lambda_n^{(n)} = \lambda_\ell + \frac{n-1}{n}(\lambda_r - \lambda_\ell)$$

we see that  $\{e^{\lambda_n^{(n)} \omega}\}$  is complete in  $L^2[\lambda_\ell, \lambda_r]$  and  $C[\lambda_\ell, \lambda_r]$ .

The following theorem due to Banach is also proven in Davis (1963):

"Let  $X$  be a normed linear space (real or complex). A sequence of elements  $\{x_k\}$  is closed if and only if it is complete."

Thus from the above theorem the sequence  $\{e^{\lambda_n^{(n)} \omega}\}$  is closed in

$L^2[\lambda_\ell, \lambda_r]$  and  $C[\lambda_\ell, \lambda_r]$ .

Finally, Davis proves the following theorem:

If  $X$  is a normed linear space,  $y$  an arbitrary element of  $X$ , and

$$E_n(y) = \min_{a_i} \left\| y - \sum_{i=1}^m a_i x_i \right\|$$

$\lim_{n \rightarrow \infty} E_n(y) = 0$  for all  $y \in X$  if and only if the independent sequence  $\{x_i\}$  is closed in  $X$ .

Thus  $\{e_n^{(n)}\}$  is dense in  $L^2[\lambda_\ell, \lambda_r]$  and in  $C[\lambda_\ell, \lambda_r]$  and hence we see that this method of estimation is tractable for the normal distribution. It is not to be understood that the particular sequence  $\{\lambda_n^{(n)}\}$  here used is the only, or in some sense best, sequence for this purpose. For example, if  $\lambda_\ell \neq -\lambda_r$ , then the sequence

$$\{\lambda_n^{(n)}\} = \lambda_\ell, \lambda_r, \frac{\lambda_\ell + \lambda_r}{2}, \frac{\lambda_\ell + \lambda_r}{4}, \frac{3(\lambda_\ell + \lambda_r)}{4}, \dots,$$

would also serve.

Now let us consider the Gamma family:

$$(2.34) \quad f(x|\lambda) = \frac{|\lambda|^r}{\Gamma(r)} x^{r-1} e^{-|\lambda|x}, \quad r > 0, x \in [0, \infty), \lambda \in (-\infty, 0).$$

For this family  $\beta(\lambda) = \frac{|\lambda|^r}{\Gamma(r)}$  and thus

$$(2.35) \quad \frac{\beta(\lambda_i) \beta(\lambda)}{\beta(\lambda_i^{(n)} + \lambda)} = \frac{|\lambda_i^{(n)}|^r |\lambda|^r}{\Gamma(r) (|\lambda_i^{(n)} + \lambda|^r)}.$$

Now consider

$$\frac{1}{|\lambda_i^{(n)} + \lambda|^r}.$$

Since  $\lambda_i^{(n)}$  and  $\lambda$  are always less than 0, let  $\lambda' = -\lambda$ , i.e.,  $\lambda' = |\lambda|$  and then  $|\lambda_i^{(n)} + \lambda| = \lambda_i^{(n)' + \lambda'$ . Now dropping the primes let

$$(2.36) \quad \lambda = (\lambda_r - \lambda_\ell) \omega + \lambda_\ell, \quad \lambda_\ell < \lambda_r,$$

so that we have to consider

$$(2.37) \quad \frac{1}{(\lambda_i^{(n)} + \lambda_\ell)^r \left[ 1 + \frac{\lambda_r - \lambda_\ell}{\lambda_i^{(n)} + \lambda_\ell} \omega \right]^r}$$

or in other words, without loss of generality

$$(2.38) \quad \frac{1}{\left[ 1 + \frac{\lambda_r - \lambda_\ell}{\lambda_i^{(n)} + \lambda_\ell} \omega \right]^r}, \quad (0 \leq \omega \leq 1).$$

Now  $F(z) = (1+z)^s$ ,  $s \neq 1, 2, \dots$ , has all Maclaurin's coefficients non-zero and the radius of convergence  $r = 1$ . Thus by Szàsz' theorem we have  $\{(1+t_n x)^s\}$  complete in  $L^2[0, 1]$  and  $C[0, 1]$  for

$$0 < |t_n| \leq 1 - \epsilon, \quad \epsilon > 0.$$

Since  $\lambda_r > \lambda_\ell$  and  $\lambda_\ell > 0$

$$(2.39) \quad \left| \frac{\lambda_r - \lambda_\ell}{\lambda_i^{(n)} + \lambda_\ell} \right| = \frac{\lambda_r - \lambda_\ell}{\lambda_i^{(n)} + \lambda_\ell}$$

and

$$(2.40) \quad \frac{\lambda_r - \lambda_\ell}{\lambda_i^{(n)} + \lambda_\ell} < 1$$

for

$$(2.41) \quad \lambda_i^{(n)} > \lambda_r - 2\lambda_\ell.$$

Thus let

$$t_n = \frac{\lambda_r - \lambda_\ell}{\lambda_n^{(n)} + \lambda_\ell}$$

where  $\lambda_n^{(n)} = n\lambda_r > \lambda_r - 2\lambda_\ell$ ,  $n = 1, 2, \dots$ , and then  $\frac{1}{(1+t_n\omega)^r}$  is complete in  $L^2[0, 1)$  and  $C[0, 1)$ , i. e.,

$$(2.42) \quad \left\{ \frac{1}{\left[ 1 + \frac{\lambda_r - \lambda_\ell}{\lambda_n^{(n)} + \lambda_\ell} \omega \right]^r} \right\} = \left\{ \frac{(\lambda_n^{(n)} + \lambda_\ell)^r}{((\lambda_n^{(n)} + \lambda_\ell) + (\lambda_r - \lambda_\ell)\omega)^r} \right\}.$$

is complete in  $L^2[0, 1]$  and  $C[0, 1]$ .

Then so is

$$\left\{ \frac{1}{((\lambda_n^{(n)} + \lambda_\ell + (\lambda_r - \lambda_\ell)\omega)^r)} \right\}$$

and thus

$$(2.43) \quad \left\{ \frac{\lambda_n^{(n)r}}{\Gamma(r)(\lambda_n^{(n)} + \lambda)^r} \right\}$$

is complete in  $L^2[\lambda_\ell, \lambda_r]$ .

Lemma. Let  $w(x) \in C[0, 1]$  and  $\frac{1}{w(x)} \geq \varepsilon > 0$  there. Then  $\{w(x)f_n(x)\}$  is closed (complete) in  $L^2[0, 1]$  if and only if  $\{f_n(x)\}$  is.

Proof. (Although the proofs of both parts of the lemma are easy we will give only the "if" part here.) Let  $g(x)$  be any element of  $L^2[0, 1]$ . Then

$$\int_0^1 \frac{g^2(x)}{w^2(x)} d\mu(x) \leq \frac{1}{\varepsilon^2} \int_0^1 g^2(x) d\mu(x) < \infty.$$

Therefore  $\frac{g(x)}{w(x)} \in L^2[0, 1]$ . Thus  $\frac{g(x)}{w(x)}$  can be approximated arbitrarily closely by finite linear combinations of elements of  $\{f_n(x)\}$ . Using these linear combinations which approximate  $\frac{g(x)}{w(x)}$  and multiplying them by  $w(x)$  we approximate  $g(x)$  by  $\{w(x)f_n(x)\}$ .

Clearly this lemma generalizes to  $L^2[\lambda_\ell, \lambda_r]$ . Let  $w(\lambda) = \lambda^r$ .

Then since

$$f_n(\lambda) = \frac{\lambda_n^r}{\Gamma(r)(\lambda_n^{(n)} + \lambda)^r}$$

we have

$$(2.44) \quad \left\{ w(\lambda) f_n(\lambda) \right\} = \left\{ \frac{\lambda_n^{(n)} \lambda^r}{\Gamma(r)(\lambda_n^{(n)} + \lambda)^r} \right\}$$

is closed (complete) in  $L^2[\lambda_\ell, \lambda_r]$ . Note that since  $\frac{1}{\lambda^r}$  is continuous on  $[\lambda_\ell, \lambda_r]$ , ( $\lambda_\ell > 0$ ) a slight modification of this proof for  $[\lambda_\ell, \lambda_r]$  yields closure in  $C[\lambda_\ell, \lambda_r]$ . Remembering that the  $\lambda$ 's above represent  $|\lambda|$  and  $\lambda_n^{(n)} + \lambda$  represents  $|\lambda_n^{(n)} + \lambda|$  for  $\lambda < 0$  we have proven the following theorem:

The  $\beta(\cdot)$  functions of the Gamma family are closed in  $L^2[\lambda_\ell, \lambda_r]$  and  $C[\lambda_\ell, \lambda_r]$ .

In accordance with the above we make the following definition:

Any member of the exponential family which is identifiable and is such that there exists a sequence

$$\left\{ \frac{\beta(\lambda_n) \beta(\lambda)}{\beta(\lambda_n + \lambda)} \right\}$$

which is closed in  $L^2[\lambda_\ell, \lambda_r]$  and  $C[\lambda_\ell, \lambda_r]$  will be called



approximable.

Let us assume now that  $G(\lambda)$  is known to be continuous.

Theorem. Let  $f(x|\lambda)$  be approximable and  $G(\lambda)$  continuous. Then given  $\varepsilon > 0$  and  $\lambda^* \in [\lambda_\ell, \lambda_r]$  there exists an  $N^*$  such that  $|E_X(\Phi_{N^*}(x|\lambda^*) - G(\lambda^*))| < \varepsilon$ .

Proof. Since  $f(x|\lambda)$  is approximable there exists a sequence

$$\left\{ \frac{\beta(\lambda_n)\beta(\lambda)}{\beta(\lambda_n + \lambda)} \right\}$$

which converges in  $L^2$ . Thus we may form a sequence  $\{H_n(\lambda|\lambda^*)\}$  such that

$$(2.45) \quad H_{n_k}(\lambda|\lambda^*) \xrightarrow{L^2} H(\lambda|\lambda^*).$$

Therefore there exists a subsequence  $\{H_{n_k}(\lambda|\lambda^*)\}$  such that

$$(2.46) \quad H_{n_k}(\lambda|\lambda^*) \xrightarrow{\text{a. e.}} H(\lambda|\lambda^*).$$

Thus we may choose a positive integer  $K$  such that

$$|H_{n_K}(\lambda|\lambda^*) - H(\lambda|\lambda^*)| < \varepsilon \quad \text{for all } \lambda \in [\lambda_\ell, \lambda_r] \setminus A$$

where  $A \subset [\lambda_\ell, \lambda_r]$  and  $m(A) = 0$ . Now by (2.11) and (2.14) we

have

$$(2.47) \quad \left| E_X(\Phi_N(x|\lambda^*)) - G(\lambda^*) \right| = \left| \int_{\lambda_\ell}^{\lambda_g} H_N(\lambda|\lambda^*) dG(\lambda) - \int_{\lambda_\ell}^{\lambda_r} H(\lambda|\lambda^*) dG(\lambda) \right|$$

where  $N = n_K$  So

$$(2.48) \quad |E_X(\Phi_N(x|\lambda^*)) - G(\lambda^*)| \leq \int_{\lambda_\ell}^{\lambda_r} |H_N(\lambda|\lambda^*) - H(\lambda|\lambda^*)| dG(\lambda).$$

Now

$$(2.49) \quad \int_{\lambda_\ell}^{\lambda_r} |H_N(\lambda|\lambda^*) - H(\lambda|\lambda^*)| dG(\lambda) \\ = \int_{[\lambda_\ell, \lambda_r] \setminus A} |H_N(\lambda|\lambda^*) - H(\lambda|\lambda^*)| dG(\lambda) + \int_A |H_N(\lambda|\lambda^*) - H(\lambda|\lambda^*)| dG(\lambda).$$

But both  $H_N(\lambda|\lambda^*)$  and  $H(\lambda|\lambda^*)$  are bounded on  $[\lambda_\ell, \lambda_r]$ . Therefore

$$(2.50) \quad \int_A |H_N(\lambda|\lambda^*) - H(\lambda|\lambda^*)| dG(\lambda) \leq M \int_A dG(\lambda)$$

where

$$|H_N(\lambda|\lambda^*) - H(\lambda|\lambda^*)| \leq M < \infty$$

for

$$\lambda \in [\lambda_\ell, \lambda_r].$$

Since  $G(\lambda)$  is continuous  $A$  is a set of  $G$  measure 0, and thus

$$(2.51) \quad M \int_A dG(\lambda) = M \cdot 0 = 0.$$

Also we have

$$(2.52) \quad \int_{[\lambda_\ell, \lambda_r] \setminus A} |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| dG(\lambda) \leq \varepsilon \int_{[\lambda_\ell, \lambda_r] \setminus A} dG(\lambda) \leq \varepsilon.$$

Therefore from (2.48), (2.49), (2.51), and (2.52) we have

$$(2.53) \quad |E_X(\Phi_N(\lambda | \lambda^*)) - G(\lambda^*)| < \varepsilon.$$

As a generalization of the above theorem suppose the points of discontinuity of  $G(\lambda)$  are known. (Actually the theorem to be proven will hold even if the discontinuity points are unknown.)

Theorem. Let  $f(x|\lambda)$  be approximable and  $\lambda^*$  be a point of continuity of  $G(\lambda)$ . Then given  $\varepsilon > 0$  there exists an  $N$  such that  $|E_X(\Phi_N(x|\lambda^*)) - G(\lambda^*)| < \varepsilon$ .

Proof. Suppose  $\lambda^*$  is a continuity point of  $G(\lambda)$ . Then given  $\varepsilon > 0$  there exists a closed interval  $[\lambda^*, \lambda^* + \delta]$  such that  $G(\lambda)$  is continuous on  $[\lambda^*, \lambda^* + \delta]$  and  $G(\lambda^* + \delta) - G(\lambda^*) < \frac{\varepsilon}{2}$ .

Consider the function

$$(2.54) \quad H^*(\lambda) = \begin{cases} 1, & \lambda_\ell \leq \lambda \leq \lambda^* \\ 1 - \frac{1}{\delta}(\lambda - \lambda^*), & \lambda^* < \lambda \leq \lambda^* + \delta \\ 0, & \lambda^* + \delta < \lambda \leq \lambda_r \end{cases}$$

Choose  $N$  so that

$$(2.55) \quad |H_N(\lambda | \lambda^*) - H^*(\lambda)| < \frac{\varepsilon}{2}, \quad \lambda \in [\lambda_\ell, \lambda_r].$$

Then

$$(2.56) \quad |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| < \frac{\varepsilon}{2}, \quad \lambda \in [\lambda_\ell, \lambda^*], \lambda \in (\lambda^* + \delta, \lambda_r]$$

and

$$(2.57) \quad |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| < 1 + \frac{\varepsilon}{2}, \quad \lambda \in (\lambda^*, \lambda^* + \delta].$$

Now from (2.48)

$$|E_X(\Phi_N(\mathbf{x} | \lambda^*)) - G(\lambda^*)| \leq \int_{\lambda_\ell}^{\lambda_r} |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| dG(\lambda).$$

Further by (2.56) and (2.57) we have

$$(2.58) \quad \int_{\lambda_\ell}^{\lambda_r} |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| dG(\lambda) \\ = \int_{\lambda_\ell}^{\lambda^*} |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| dG(\lambda) + \int_{\lambda_\ell}^{\lambda^* + \delta} |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| dG(\lambda) +$$

$$\begin{aligned}
& + \int_{\lambda^* + \delta}^{\lambda^r} |H_N(\lambda | \lambda^*) - H(\lambda | \lambda^*)| dG(\lambda) \\
& \leq \frac{\varepsilon}{2} \int_{\lambda_\ell}^{\lambda^*} dG(\lambda) + \frac{\varepsilon}{2} \int_{\lambda^* + \delta}^{\lambda^r} dG(\lambda) + (1 + \frac{\varepsilon}{2}) \int_{\lambda^*}^{\lambda^* + \delta} dG(\lambda) \\
& = \frac{\varepsilon}{2} \int_{\lambda_\ell}^{\lambda^r} dG(\lambda) + \int_{\lambda^*}^{\lambda^* + \delta} dG(\lambda) \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon.
\end{aligned}$$

Thus by (2.48) and (2.58)

$$(2.59) \quad |E_X(\Phi_N(x | \lambda^*)) - G(\lambda^*)| \leq \varepsilon.$$

Since  $\lambda^r$ ,  $r \in I^+$ , is continuous on  $[\lambda_\ell, \lambda_r]$  we may approximate any  $r^{\text{th}}$  moment of  $\lambda$ .

We have the

Theorem. Let  $f(x|\lambda)$  be approximable. Then given  $\varepsilon > 0$  there exists an  $N$  such that  $|E_X(\Phi_N(x|\lambda_r)) - E_\Lambda(\lambda^r)| < \varepsilon$ .

Proof. In (2.14) replace  $\lambda^*$  by  $\lambda_r$  and let

$$(2.60) \quad H(\lambda | \lambda_r) = \lambda^r.$$

Thus in (2.17) we have

$$(2.61) \quad H_n(\lambda | \lambda_r) = \sum_{i=1}^{k_n} a_{i\lambda_r}^{(n)} h_{\lambda_i}^{(n)}(\lambda).$$

Now since  $f(x|\lambda)$  is approximable there exists an  $N$  such that

$$(2.62) \quad |H_N(\lambda | \lambda_r) - H(\lambda | \lambda_r)| < \varepsilon$$

we have

$$(2.63) \quad |E_X(\Phi_N(x|\lambda_r)) - E_\Lambda(\lambda^r)| = \left| \int_\Lambda H_N(\lambda | \lambda_r) dG(\lambda) - \int_\Lambda \lambda^r dG(\lambda) \right|$$

$$\leq \int_\Lambda |H_N(\lambda | \lambda_r) - H(\lambda | \lambda_r)| dG(\lambda)$$

$$\leq \varepsilon \int_\Lambda dG(\lambda) = \varepsilon.$$

From what has been done above it is clear that we may replace  $H(\lambda | \lambda_r^*)$  by any function of  $\lambda$ ,  $\psi(\lambda)$ , which is continuous on  $[\lambda_\ell, \lambda_r]$  and thus estimate  $E_\Lambda[\psi(\lambda)]$ .

At (2.8) we defined

$$\bar{G}_m^n(\lambda^*) = \frac{1}{m} \sum_{v=1}^m \Phi_n(x_v | \lambda^*).$$

This function is not necessarily a distribution function since it may not be monotonic increasing, nor is it necessarily non-negative

nor bounded by one. Thus it is necessary to define a function  $G_m^n(\lambda^*)$  which will be a distribution function.

Let  $C_G =$  set of continuity points of  $G(\lambda)$  in  $[\lambda_\ell, \lambda_r]$ . We have seen by a previous theorem that given  $\varepsilon > 0$ ,  $\lambda^* \in C_G$  there exists an  $n(\varepsilon, \lambda^*)$  such that  $|E_X(\Phi_{n(\varepsilon, \lambda^*)}(x|\lambda^*)) - G(\lambda^*)| < \varepsilon$ .

Lemma. Given  $\varepsilon > 0$  let

$$(2.64) \quad N_\varepsilon = \sup_{\lambda^* \in C_G} \{n(\varepsilon, \lambda^*)\}.$$

Then  $N_\varepsilon < \infty$ .

Proof. Suppose not. Then there exists  $\lambda^* \in C_G$  such that

$$|E_X(\Phi_n(x|\lambda^*)) - G(\lambda^*)| \geq \varepsilon$$

for all  $n \in I^+$ . But this contradicts the completeness of

$$\left\{ \frac{\beta(\lambda_n)\beta(\lambda)}{\beta(\lambda + \lambda_n)} \right\}$$

in  $C[\lambda_\ell, \lambda_r]$ . Thus  $N_\varepsilon < \infty$ . We may write

$$(2.65) \quad \bar{G}_m^{N_\varepsilon}(\lambda^*) = \frac{1}{m} \sum_{\nu=1}^m \Phi_{N_\varepsilon}^{(\nu)}(x_\nu|\lambda^*),$$

and

$$|E_X(\Phi_{N_\varepsilon}(x|\lambda^*)) - G(\lambda^*)| < \varepsilon$$

for all  $\lambda^* \in C_G$ . Now let

$$(2.66) \quad E_X(\Phi_{N_\varepsilon}(x|\lambda^*)) = F_\varepsilon^N(\lambda^*).$$

Then we have the

Theorem.

$$\bar{G}_m^N(\lambda^*) \xrightarrow{\text{a. s.}} F_\varepsilon^N(\lambda^*)$$

where

$$G(\lambda^*) - \varepsilon \leq F_\varepsilon^N(\lambda^*) \leq G(\lambda^*) + \varepsilon$$

for all  $\lambda^* \in C_G$ .

Proof. By the Kolmogorov strong law of large numbers

$$\bar{G}_m^N(\lambda^*) \xrightarrow{\text{a. s.}} F_\varepsilon^N(\lambda^*)$$

and by the preceding theorem

$$G(\lambda^*) - \varepsilon \leq F_\varepsilon^N(\lambda^*) \leq G(\lambda^*) + \varepsilon.$$

Let



$$(2.67) \quad [\bar{G}_m^N(\lambda^*)] = \begin{cases} 0, & \bar{G}_m^N(\lambda^*) \leq 0 \\ \bar{G}_m^N(\lambda^*), & 0 < \bar{G}_m^N(\lambda^*) < 1 \\ 1, & \bar{G}_m^N(\lambda^*) \geq 1 \end{cases}$$

and

$$G_m^N(\lambda^*) = \begin{cases} 0, & \lambda^* < \lambda_l \\ \max_{\lambda_l \leq \theta \leq \lambda^*} [\bar{G}_m^N(\theta)], & \lambda_l \leq \lambda^* \leq \lambda_r \\ 1, & \lambda^* > \lambda_r \end{cases}$$

Clearly  $G_m^N(\lambda^*) \xrightarrow{\text{a.s.}} F^N(\lambda^*)$  and further  $G_m^N(\lambda^*)$  is a distribution function. Thus we have proved the following

Theorem.  $G_m^N(\lambda^*)$  is a distribution function such that

(2.69)

$$P \left[ \lim_{m \rightarrow \infty} G_m^N(\lambda^*) = F^N(\lambda^*) \text{ at every continuity point of } G(\lambda^*) \right] = 1,$$

where

$$G(\lambda^*) - \varepsilon \leq F^N(\lambda^*) \leq G(\lambda^*) + \varepsilon.$$

### III. ESTIMATION OF DISCRETE PRIOR DISTRIBUTIONS WITH VARYING DIMENSIONAL RANDOM VECTORS

Consider the measure space  $(X, \mathcal{X}, \mu)$  with  $\mu$   $\sigma$ -finite and the random variable  $x \in X$  known to have one of a finite number of probability distributions

$$(3.1) \quad P_i, \quad 1 \leq i \leq r,$$

such that all  $P_i$  are absolutely continuous with respect to  $\mu$  and their densities  $f_i = \frac{dP_i}{d\mu}$  square integrable and bounded a. e.  $\mu(x)$ .

Which  $P_i$  the random variable  $x$  has depends on a random parameter  $\lambda \in \Lambda = \{\lambda_1, \dots, \lambda_r\}$  which has an unknown a-priori probability vector

$$(3.2) \quad G = \{g_1, g_2, \dots, g_r\},$$

$$g_i \geq 0, \quad \sum_{i=1}^r g_i = 1$$

such that

$$(3.3) \quad P(\lambda = \lambda_i) = g_i.$$

We observe a sequence of random vectors  $x_1^{(v_1)}, x_2^{(v_2)}, \dots$

as follows: the values of  $\lambda_j$  are drawn independently from  $G(\lambda)$

and the random sample sizes  $v_j$  and determined by the sequence

$\{x_k^{(v_k)}\}_{k=1}^{j-1}$ . Given  $\lambda_j$  and  $v_j$ , the  $v_j$  components of  $x_j^{(v_j)}$

are conditionally independent with common density function  $f(x|\lambda_j)$ , with

$$(3.4) \quad x_j^{(v_j)} \in (X^{(v_j)}, \chi^{(v_j)}, \mu^{(v_j)})$$

the  $v_j$  fold product space generated by  $v_j$  fold Cartesian products of  $(X, \chi, \mu)$ ,  $1 \leq v_j \leq N$ . That is, at the  $j$ th instant we have a choice of observing any one of  $N$  random vectors

$$(3.5) \quad x_j^{(v_j=1)} \in X, \quad x_j^{(v_j=2)} \in X \times X, \dots, \\ x_j^{(v_j)} \in \underbrace{X \times X \times \dots \times X}_{v_j \text{ fold}}, \dots, x_j^{(v_j=N)} \in \underbrace{X \times \dots \times X}_{N \text{ fold}}$$

with the dimension  $v_j$  of the vector observed at the  $j$ th instant being determined by a rule depending on the preceding observations. The  $j$ th of these random vectors has the conditional distribution, given  $v_j$ ,

$$(3.6) \quad P_G^{(v_j)}(x_j^{(v_j)} \in B^{(v_j)}) = \sum_{i=1}^r g_i P_i^{(v_j)}(B^{(v_j)}),$$

where  $B^{(v_j)} \in \chi^{(v_j)}$  and  $P_i^{(v_j)}(\cdot)$  is the probability distribution generated on  $(X^{(v_j)}, \chi^{(v_j)}, \mu^{(v_j)})$  by the  $P_i(\cdot)$  on  $(X, \chi, \mu)$ .

Our problem is to construct functions

$$(3.7) \quad g_{i, n} = g_{i, n}^{(v_1)}(x_1), x_2^{(v_2)}, \dots, x_n^{(v_n)}$$

such that

$$g_{i, n} \geq 0, \quad \sum_{i=1}^r g_{i, n} = 1,$$

and whatever be  $G$ ,

$$(3.8) \quad P \left[ \lim_{n \rightarrow \infty} g_{i, n} = g_i \right] = 1, \quad (i = 1, 2, \dots, r).$$

Since we are dealing with random samples the components of the random vectors are independent random variables and the conditional (given  $v_j$ ) density function of  $x_i^{(v_j)}$  is then

$$(3.9) \quad f_i^{(v_j)}(x_j^{(v_j)}) = \prod_{k=1}^{v_j} f_i(x_k).$$

Now due to (3.9) and the fact that

$$\int_X f_i^2(x) d\mu(x) < \infty$$

we have

$$(3.10) \quad \int_{X^{(v_j)}} [f_i^{(v_j)}(x_j^{(v_j)})]^2 d\mu^{(v_j)}(x_j^{(v_j)}) < \infty, \quad (i = 1, 2, \dots, r)$$

Thus the functions  $f_i^{(v_j)}$  are elements of the Hilbert space  $H^{(v_j)}$  of square integrable functions over the measure space  $(X^{(v_j)}, \mathcal{X}^{(v_j)}, \mu^{(v_j)})$ .

We make at this point the further assumption of identifiability, i. e., if  $G = \{g_1, \dots, g_r\}$  and  $\bar{G} = \{\bar{g}_1, \dots, \bar{g}_r\}$  are any two probability vectors such that for every  $B_j^{(v_j)}$ ,

$$\sum_{i=1}^r g_i P_i^{(v_j)}(B_j^{(v_j)}) = \sum_{i=1}^r \bar{g}_i P_i^{(v_j)}(B_j^{(v_j)}),$$

then  $G = \bar{G}$ . The class of measures  $P$  which are identifiable under such finite mixtures is a subclass of those which are identifiable under general mixtures.

Lemma 1 of Chapter II allows us to denote by  $H_k^{(v_j)}$  the linear manifold spanned by the  $r - 1$  functions  $f_1^{(v_j)}, \dots, f_{k-1}^{(v_j)}, f_{k+1}^{(v_j)}, \dots, f_r^{(v_j)}$ . We can then write uniquely

$$(3.11) \quad f_k^{(v_j)} = f_k^{(v_j)'} + f_k^{(v_j)''}, \quad (k = 1, 2, \dots, r)$$

where

$$(3.12) \quad f_k^{(v_j)'} \in H_k^{(v_j)}, \quad f_k^{(v_j)''} \in H_k^{(v_j)\perp}, \quad \text{and} \quad f_k^{(v_j)''} \neq 0.$$

Let

$$(3.13) \quad \phi_{kj}(x^{(v_j)}) = \frac{f_k^{(v_j)''}(x^{(v_j)})}{\int_{X^{(v_j)}} [f_k^{(v_j)''}(x^{(v_j)})]^2 d\mu^{(v_j)}(x^{(v_j)})},$$

then

$$(3.14) \quad \int_{\mathbf{X}} \phi_{kj}^{(v_j)}(t) f_{\ell}^{(v_j)}(t) d\mu = \begin{cases} 1, & \text{if } k = \ell, \\ 0, & \text{if } k \neq \ell \end{cases}.$$

Now define

$$(3.15) \quad \bar{g}_{i,n} = \frac{1}{n} \sum_{j=1}^n \phi_{ij}^{(v_j)}(x_j),$$

$$g_{i,n} = \frac{[\bar{g}_{i,n}]^+}{\sum_{j=1}^r [\bar{g}_{i,n}]^+} \cdot ([a]^+ = \max(a, 0))$$

Note that any convergence properties of  $\bar{g}_{i,n}$  will be shared by  $g_{i,n}$ . Now

$$(3.16) \quad E_{\mathbf{X}} \phi_{ij}^{(v_j)}(x_j) = \int_{\mathbf{X}} \phi_{ij}^{(v_j)}(x_j) \sum_{k=1}^r g_k f_k^{(v_j)}(x_j) d\mu^{(v_j)}(x_j)$$

$$= \sum_{k=1}^r g_k \int_{\mathbf{X}} \phi_{ij}^{(v_j)}(x_j) f_k^{(v_j)}(x_j) d\mu^{(v_j)}(x_j)$$

$$= g_i.$$

Note that whatever value of the random variable  $v_j$  may occur the conditional expectation of  $\phi_{ij}^{(v_j)}(x_j)$  is  $g_i$ . Thus the unconditional expectation of  $\phi_{ij}^{(v_j)}(x_j)$  is  $g_i$ . Also

$$\begin{aligned}
 (3.17) \quad E_{X^{(v_j)}} \phi_{ij}^2(x_j^{(v_j)}) &= \int_{X^{(v_j)}} \phi_{ij}^2(x_j^{(v_j)}) \sum_{k=1}^r g_k f_k^{(v_j)}(x_j^{(v_j)}) d\mu^{(v_j)}(x_j^{(v_j)}) \\
 &= \sum_{k=1}^r g_k \int_{X^{(v_j)}} \phi_{ij}^2(x_j^{(v_j)}) f_k^{(v_j)}(x_j^{(v_j)}) d\mu^{(v_j)}(x_j^{(v_j)}).
 \end{aligned}$$

By the initially assumed conditions there exists  $M_k < \infty$  such that

$$(3.18) \quad 0 \leq f_k^{(v_j)}(x_j^{(v_j)}) \leq M_k < \infty, \quad \text{a. e. } \mu^{(v_j)}(x_j^{(v_j)}),$$

therefore

$$(3.19) \quad E_{X^{(v_j)}} \phi_{ij}^2(x_j^{(v_j)}) \leq \sum_{k=1}^r g_k M_k \int_{X^{(v_j)}} \phi_{ij}^2 d\mu^{(v_j)}(x_j^{(v_j)}).$$

But

$$\begin{aligned}
 (3.20) \quad &\int_{X^{(v_j)}} \phi_{ij}^2 d\mu^{(v_j)}(x_j^{(v_j)}) \\
 &= \frac{\int_{X^{(v_j)}} \left[ f_i^{(v_j)}(x_j^{(v_j)}) \right]^2 d\mu^{(v_j)}(x_j^{(v_j)})}{\int_{X^{(v_j)}} \left[ f_i^{(v_j)}(x_j^{(v_j)}) \right]^2 d\mu^{(v_j)}(x_j^{(v_j)})} \\
 &= \frac{1}{\int_{X^{(v_j)}} \left[ f_i^{(v_j)}(x_j^{(v_j)}) \right]^2 d\mu^{(v_j)}(x_j^{(v_j)})},
 \end{aligned}$$

and from (3.12), we have

$$(3.21) \quad \int_{X_j}^{(\nu_j)} \left[ f_i^{(\nu_j)}(x_j) \right]^2 d\mu_j^{(\nu_j)}(x_j) > 0,$$

so that

$$(3.22) \quad \frac{1}{\int_{X_j}^{(\nu_j)} \left[ f_i^{(\nu_j)}(x_j) \right]^2 d\mu_j^{(\nu_j)}(x_j)} = M_{ij}^* < \infty.$$

This produces

$$(3.23) \quad E_{X_j}^{(\nu_j)} \phi_{ij}^2(x_j) \leq \sum_{k=1}^r g_k M_k M_{ij}^*$$

and letting  $M' = \max_k \{M_k\}$  and  $M^* = \max_{i,j} \{M_{ij}^*\}$  we have

$$(3.24) \quad E_{X_j}^{(\nu_j)} \phi_{ij}^2(x_j) \leq M' M^* < \infty.$$

Now

$$\text{Var } \phi_{ij}(x_j) \leq E_{X_j}^{(\nu_j)} \phi_{ij}^2(x_j) < \infty.$$

So

$$(3.25) \quad \sum_{j=1}^{\infty} \frac{\text{Var } \phi_{ij}(x_j)}{j^2} \leq \sum_{j=1}^{\infty} E \phi_{ij}^2(x_j) \leq M' M^* \sum_{j=1}^{\infty} \frac{1}{j^2} < \infty.$$



Feller (1966, p. 238) gives the theorem: Let  $\{x_k\}$  be a sequence of random variables with  $E(x_k | x_1, \dots, x_{k-1}) = 0$ , and let

$$S_n = \sum_{k=1}^n x_k.$$

If  $b_1 < b_2 < \dots \rightarrow \infty$  and

$$(3.26) \quad \sum_{k=1}^{\infty} E(x_k^2)/b_k^2 < \infty,$$

Then

$$S_n/n \xrightarrow{\text{a.s.}} 0.$$

If we let  $x_k$  represent  $\phi_{ij}(x_j^{(v_j)}) - g_i$  and  $b_n = n$ , then (3.25) implies that (3.26) is satisfied. Furthermore (3.16) implies that

$E(x_k | x_1, \dots, x_{k-1}) = 0$ , therefore we have

$$(3.27) \quad \frac{1}{n} \sum_{j=1}^n \phi_{ij}(x_j^{(v_j)}) - g_i \xrightarrow{\text{a.s.}} 0$$

and

$$(3.28) \quad \frac{1}{n} \sum_{j=1}^n \phi_{ij}(x_j^{(v_j)}) \xrightarrow{\text{a.s.}} g_i.$$

We have therefore completed the construction of consistent estimates

of the  $g_i$ .

Guthrie and Johns (1959) have studied the problem of choice of sample size in a Bayesian approach to lot-by-lot acceptance sampling. They derive asymptotic expressions for the optimal sample sizes and decision procedures as the lot size becomes large, assuming throughout that the a-priori distribution of expected lot quality is known. They study two cases: when the a-priori distribution is continuously differentiable with respect to Lebesgue measure; and when the a-priori distribution belongs to a family of atomic distributions whose support contains no point of accumulation at a certain critical point. Included in the latter class is the class of distributions over a finite number of points.

A critical parameter in the development of the asymptotic sample size is

$$A_G = (s_1 - r_1)E(\Lambda) + s_2 - r_2 + (r_1 - a_1) \int_0^C (\lambda - C) dG(\lambda)$$

where  $r_1$ ,  $r_2$ ,  $s_1$ ,  $s_2$ , and  $C$  are given known constants.

We may use the previous development to give empirical estimates for  $A_G$  in the case where many lots are inspected, one at a time, with each sample size based on the previously estimated value of  $A_G$ . Let

$$(3.29) \quad E_n(\lambda) = \sum_{i=1}^r \lambda_i g_{i,n}.$$

Then clearly

$$(3.30) \quad E_n(\lambda) \xrightarrow{\text{a. s.}} E(\lambda).$$

Also

$$(3.31) \quad \int_0^C (\lambda - C) dG(\lambda) = \sum_{i=1}^{[k]} (\lambda_i - C)(g_i - g_{i-1})$$

where  $[k]$  = the greatest integer such that  $\lambda_i \leq C$ . So

$$(3.32) \quad \sum_{i=1}^{[k]} (\lambda_i - C)(g_{i,n} - g_{i-1,n}) \xrightarrow{\text{a. s.}} \sum_{i=1}^{[k]} (\lambda_i - C)(g_i - g_{i-1}) = \int_0^C (\lambda - C) dG(\lambda).$$

Thus we define

$$A_{G_n} = (s_1 - r_1)E_n(\lambda) + s_2 - r_2 + (r_1 - a_1) \sum_{i=1}^{[k]} (\lambda_i - C)(g_{i,n} - g_{i-1,n})$$

and we have

$$(3.33) \quad A_{G_n} \xrightarrow{\text{a. s.}} A_G.$$

If we will now let

$$(3.34) \quad v_n = \begin{cases} N, & A_{G_n} \leq 0 \\ K_n \ln N - \frac{K_n}{2} \ln \ln N, & A_{G_n} > 0 \end{cases}$$

(For definition of  $K_n$  see definition of  $K$  on p. 922 of the Guthrie-Johns paper), we immediately see that  $\nu_n \xrightarrow{\text{a.s.}} \nu = \nu^*(N)$ , the approximate Bayes sample size for a-priori  $G$ , omitting the remainder terms in the Guthrie-Johns expansion.

It is readily verified that their loss structure satisfies

$$(3.35) \quad 0 \leq L(a, \lambda_i) \leq L < \infty, \quad \forall a \in A, \quad i = 1, 2, \dots, r.$$

We therefore have

$$(3.36) \quad \Delta_G(a, x^{(\nu)}) = \sum_{i=1}^r [L(a, \lambda_i) - L(a_0, \lambda_i)] \cdot f_i^{(\nu)}(x^{(\nu)}) g_i.$$

Set

$$(3.37) \quad \Delta_n(a, x^{(\nu_n)}) = \sum_{i=1}^r [L(a, \lambda_i) - L(a_0, \lambda_i)] \cdot f_i^{(\nu_n)}(x^{(\nu_n)}) g_{i,n},$$

so that

$$(3.38) \quad \sup_A |\Delta_n(a, x^{(\nu_n)}) - \Delta_G(a, x^{(\nu)})| \leq L \sum_{i=1}^r |f_i^{(\nu_n)}(x^{(\nu_n)}) g_{i,n} - f_i^{(\nu)}(x^{(\nu)}) g_i|.$$

Since  $\nu_n \rightarrow \nu$ ,  $f_i^{(\nu_n)}(x^{(\nu_n)})$  and  $f_i^{(\nu)}(x^{(\nu)})$  are  $< \infty$  a.e., we see that

$$(3.39) \quad p \lim_{n \rightarrow \infty} \sup_A |\Delta_n(a, x^{(\nu_n)}) - \Delta_G(a, x^{(\nu)})| = 0$$

Thus our sequence of decision functions is asymptotically optimal in the sense of Robbins, hence the empirical acceptance sampling procedure so generated is asymptotically equivalent to that for known a-priori distributions, omitting remainder terms.

## BIBLIOGRAPHY

- Cheney, E. W. 1966. Introduction to approximation theory. New York, McGraw-Hill. 259 p.
- Davis, P. J. 1963. Interpolation and approximation. New York, Blaisdell. 385 p.
- Feller, William. 1966. An introduction to probability theory and its applications. Vol. 2. New York, Wiley. 626 p.
- Guthrie, D., Jr. and M. V. Johns, Jr. 1959. Bayes acceptance sampling procedures for large lots. *Annals of Mathematical Statistics* 30:896-925.
- Johns, M. V., Jr. 1957. Non-parametric empirical Bayes procedures. *Annals of Mathematical Statistics* 28:649-669.
- Lehmann, E. L. 1959. Testing statistical hypotheses. New York, Wiley. 369 p.
- Loeve, M. 1960. Probability theory. 2d ed. Princeton, Van Nostrand. 685 p.
- Robbins, Herbert. 1956. An empirical Bayes approach to statistics. In: *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Berkeley, 1955*. Vol. 1. Berkeley, University of California. p. 157-163.
- \_\_\_\_\_ 1964. The empirical Bayes approach to statistical decision problems. *Annals of Mathematical Statistics* 35:1-20.
- Rolph, John E. 1968. Bayesian estimation of mixing distributions. *Annals of Mathematical Statistics* 39:1289-1302.
- Samuel, Esther. 1963. An empirical Bayes approach to the testing of certain parametric hypotheses. *Annals of Mathematical Statistics* 34:1370-1385.
- Teicher, Henry. 1960. On the mixtures of distributions. *Annals of Mathematical Statistics* 31:55-73.
- \_\_\_\_\_ 1961. Identifiability of mixtures. *Annals of Mathematical Statistics* 32:244-248.

---

1963. Identifiability of finite mixtures. *Annals of Mathematical Statistics* 34:1265-1269.

Tucker, Howard G. 1967. *A graduate course in probability*. New York, Academic. 273 p.

Waterman, Michael Spencer. 1966. *The exponential family of probability distributions generated by  $\sigma$ -infinite measures*. Master's thesis. Corvallis, Oregon State University. 50 numb. leaves.