

## MEASURES OF INFORMATION BASED ON COMPARISON WITH TOTAL INFORMATION AND WITH TOTAL IGNORANCE

BY ERIK N. TORGERSEN

*University of California, Berkeley and University of Oslo*

Two measures of content of information in statistical experiments are considered. They are both based on Le Cam's notion of deficiency of one experiment with respect to another.

The measures are:

- (i) The deficiency of the given experiment with respect to a totally informative experiment.
- (ii) The deficiency of a totally uninformative experiment with respect to the given experiment.

We shall here discuss the interpretations of such measures, establish inequalities for them and related quantities, and study their behaviour under replications.

If the parameter set is finite then closed expressions for exponential rates of convergence, as the number of replications increase, are given. In particular the exponential rate of the minimax probability of not covering the true values of the parameter by an  $r$ -point confidence set is expressed in terms of Hellinger transforms. If convergence to the totally informative experiment takes place at all, then the speed of convergence is necessary exponential. Examples are given indicating various possibilities.

**1. Introduction.** In this paper we shall consider two measures of the contents of information in statistical experiments. They are both based on Le Cam's [14] notion of a deficiency of one experiment with respect to another.

Let us agree, following Le Cam [14], to use the notation  $\delta(\mathcal{E}, \mathcal{F})$  for the deficiency of an experiment  $\mathcal{E}$  with respect to an experiment  $\mathcal{F}$ . This deficiency is defined for pairs  $(\mathcal{E}, \mathcal{F})$  of experiments having the same parameter set  $\Theta$ .  $\delta(\mathcal{E}, \mathcal{F})$  is a function of two variables  $\mathcal{E}$  and  $\mathcal{F}$ . It provides a partial answer to the question: What do we lose, in the sense of risk, by basing ourselves on  $\mathcal{E}$  rather than on  $\mathcal{F}$  under the least favorable conditions for this comparison?

The deficiency  $\delta(\mathcal{E}, \mathcal{F})$  is monotonically increasing in  $\mathcal{F}$  and monotonically decreasing in  $\mathcal{E}$ . It is also convex in each variable separately. Convexity is then defined in terms of mixtures of experiments [30].

Deficiencies, or the related distances, may be considered as measures of contents of information in many situations. It may be used, see Lindqvist [19], to measure the loss of memory of the initial State  $X_0$  in the tail  $(X_t, X_{t+1}, \dots)$  of a Markov chain  $(X_0, X_1, \dots)$ . Swensen [22] has investigated the problem of measuring the value of a potential additional variable in a regression model. It is also possible to construct local measures of information based on deficiencies [24] but we shall not dwell on this here.

In order to investigate the properties of such measures it is tempting to consider, in spite of their artificiality, distances to experiments which are either totally informative or totally uninformative. A totally informative experiment is an experiment  $(P_\theta; \theta \in \Theta)$  such that  $P_{\theta_1}$  is  $P_{\theta_2}$  singular when  $\theta_1 \neq \theta_2$ . As any two totally informative experiments are

---

Received June, 1977; revised April, 1978; July, 1980.

AMS 1970 subject classifications. Primary 62B15; secondary 62C05.

Key words and phrases. Risk, replications, exponential convergence, confidence sets, translation experiments.

equivalent, we shall use the symbol  $\mathcal{M}_a$  to denote any of them. To fix ideas we might, if we so prefer, let  $\mathcal{M}_a$  denote the experiment  $(\delta_\theta; \theta \in \Theta)$  where  $\delta_\theta$  is the one point distribution in  $\theta$ . Intuitively  $\mathcal{M}_a$  is the experiment consisting of observing the underlying theory  $\theta$  itself.

A totally uninformative experiment is an experiment  $(P_\theta; \theta \in \Theta)$  where  $P_\theta$  does not depend on  $\theta$ . Clearly any two noninformative experiments are also equivalent and we shall reserve the notation  $\mathcal{M}_i$  for any of them.

Any experiment  $\mathcal{E}$  is obviously at least as informative as  $\mathcal{M}_i$  and at most as informative as  $\mathcal{M}_a$ .

In this paper we shall consider the numbers  $\delta(\mathcal{M}_i, \mathcal{E})$  and  $\delta(\mathcal{E}, \mathcal{M}_a)$  as measures of the content of information in an experiment  $\mathcal{E}$ . The deficiency  $\delta(\mathcal{M}_i, \mathcal{E})$  will usually be written  $\delta_i(\mathcal{E})$  while the deficiency  $\delta(\mathcal{E}, \mathcal{M}_a)$  will usually be written  $\delta_a(\mathcal{E})$ .

A small value of  $\delta_a(\mathcal{E})$  suggests that an observation of  $\mathcal{E}$ , provided that it is properly used, is almost as good as knowing the unknown parameter. A large value, on the other hand, tells us that there are decision problems such that any decision procedure is risky for some of the underlying theories. A small value of  $\delta_i(\mathcal{E})$  tells us that the chance mechanism governing the random outcome is almost independent of the various explanatory theories in  $\Theta$ . If, on the other hand, this distance is large, then there are situations where an observation of  $\mathcal{E}$  is helpful.

The values of these deficiencies are often large for all experiments  $\mathcal{E}$  under consideration. This reflects the fact that it may be much too ambitious to compare with total information and much too modest to compare with no information.

Having made this limitation, we shall see that these deficiencies have interesting properties.

The contents of the paper are as follows. A few of the basic definitions and results from the theory of comparison of experiments are summarized in Section 2. It is shown in Section 2 that  $\frac{1}{2} \delta(\mathcal{E}, \mathcal{M}_a)$  is the minimax probability of an incorrect guess of the true value of  $\theta$ . Similarly the deficiency  $\delta(\mathcal{M}_i, \mathcal{E})$  is the minimax risk for the problem of guessing the true value of  $\theta$ , when no observations are available and the loss is measured by statistical distance. If we restrict attention to testing problems then the corresponding deficiency reduces to the half diameter of  $\mathcal{E}$  for statistical distance.

The case of dichotomies, i.e., the case where  $\theta$  has two elements, is investigated in Section 3. We begin by slightly completing previously known comparison criteria. In the case of a finite sample space and one distribution being uniform, the results generalize some of the basic inequalities from the theory of majorization.

The remainder of the paper is devoted to replicated experiments. Let  $\mathcal{E}^n$  denote the experiment obtained by combining  $n$  independent replications of  $\mathcal{E}$ . All limits, if not otherwise stated, are taken as  $n \rightarrow \infty$ . How do these quantities behave under replications? As is well known, it follows from the weak law of large numbers that  $\mathcal{E}^n \rightarrow \mathcal{M}_a$  provided that  $P_{\theta_1} \neq P_{\theta_2}$  when  $\theta_1 \neq \theta_2$  and that  $\Theta$  is finite. If  $\Theta$  is infinite, however, then one is tempted to conclude that "normally"  $\delta_a(\mathcal{E}^n) \equiv_n 2$ . There are nevertheless interesting experiments with  $\Theta$  infinite where  $\delta_a(\mathcal{E}^n) \rightarrow 0$ .

We shall see in Section 5 that  $\delta_a(\mathcal{E}^n)^{1/n}$  converges, for any experiment  $\mathcal{E}$ , to a constant  $\sigma(\mathcal{E})$  in  $[0, 1]$ . This implies that if  $\delta_a(\mathcal{E}^n) \rightarrow 0$  then the speed of convergence is exponential. The rate of convergence is not determined by  $\delta_a(\mathcal{E})$  alone, since there are experiments  $\mathcal{E}$  such that  $\delta_a(\mathcal{E})$  has the maximal value 2 although  $\delta_a(\mathcal{E}^n)$  converges rapidly to zero.

Now Chernoff [6] proved, when  $\Theta = \{1, 2\}$ , that the  $n$ th root of the minimum Bayes probability of error converges to the minimum of the Hellinger transform. It follows that  $\sigma(\mathcal{E})$  and this minimum are the same number when  $\mathcal{E}$  is a dichotomy.

Because pairwise equivalence for ordered experiments imply equivalence, one might hope that comparison with respect to  $\mathcal{M}_i$  and  $\mathcal{M}_a$  may, to some extent, be expressed in terms of the dichotomies defined by restrictions to pairs. If  $\Theta$  is infinite, then — as shown in Section 5 — this does not hold in general. If  $\Theta$  is finite, however, then — as is shown in Section 4 — the approximations are readily expressed in terms of dichotomies. As in

Section 3 we also get the exponential rate of convergence of many other functionals.  $\sigma(\mathcal{E})$  does not, however, define the minimax probability of not covering the true value of  $\theta$  with a  $r$ -point confidence set when  $r \geq 2$ . Generalizing Chernoff's result we obtain the exponential rates of convergence of these minimax probabilities (and of the corresponding minimum Bayes probabilities).

The paper concludes with an example showing that dramatic improvement may be obtained by adding a single replication. This brings us close to a related topic — the relative amount of information in additional observations. We refer the reader to Helgeland [12], Le Cam [15], Swensen [22] and Torgersen [25] for some results in this direction.

**2. Notations and basic facts.** An experiment  $\mathcal{E}$  will be defined as a family of probability measures on a common measurable space. This measurable space and the index set of the family are called, respectively, *the sample space of  $\mathcal{E}$*  and *the parameter set of  $\mathcal{E}$* . Thus an experiment  $\mathcal{E}$  with sample space  $(\chi, \mathcal{A})$  and parameter set  $\Theta$  is a family  $(P_\theta; \theta \in \Theta)$  of probability measures on  $(\chi, \mathcal{A})$ . This experiment may be denoted by  $(P_\theta; \theta \in \Theta)$  or by  $(\chi, \mathcal{A}; P_\theta; \theta \in \Theta)$ . Expositions of the theory of the comparison of experiments may be found in Blackwell and Girshick [4], Heyer [13], Le Cam [17] and Torgersen [28].

New experiments may be derived from old ones by various devices. If  $\Theta_0$  is a subset of  $\Theta$  and  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  is an experiment with parameter set  $\Theta$ , then  $\mathcal{E}_{\Theta_0}$  denotes the restriction  $(P_\theta; \theta \in \Theta_0)$  of  $\mathcal{E}$  to  $\Theta_0$ . The restrictions of  $\mathcal{M}_a$  and  $\mathcal{M}_i$  to  $\Theta_0$  will, however, usually be denoted by  $\mathcal{M}_a$  and  $\mathcal{M}_i$ , respectively.

If  $\mathcal{E}_i = (\chi_i, \mathcal{A}_i, P_{\theta_i}; \theta \in \Theta_i); 1 \leq i \leq n$  are experiments, then  $(\prod_{i=1}^n (\chi_i, \mathcal{A}_i), \prod_{i=1}^n P_{\theta_i}; \theta \in \Theta)$  is called the *product* of  $\mathcal{E}_1, \mathcal{E}_2, \dots, \mathcal{E}_n$  and it is denoted  $\mathcal{E}_1 \times \dots \times \mathcal{E}_n$  or  $\prod \mathcal{E}_i$ . If  $\mathcal{E}_1 = \mathcal{E}_2 = \dots = \mathcal{E}_n = \mathcal{E}$  then we may write  $\mathcal{E}^n$  instead of  $\mathcal{E}_1 \times \dots \times \mathcal{E}_n$ . Experiments  $\mathcal{E}^n, n = 1, 2, \dots$  are called *replicated experiments*.

Let  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  be an experiment. A random variable  $X$  will be called an *observation of  $\mathcal{E}$*  if the distribution of  $X$  under  $\theta$  is  $P_\theta$ . Note that  $(X_1, X_2, \dots, X_n)$  is an observation of  $\mathcal{E}_1 \times \dots \times \mathcal{E}_n$  if  $X_1, \dots, X_n$  are independent and  $X_i; i = 1, \dots, n$  is an observation of  $\mathcal{E}_i$ .

If  $\Pi$  is a probability distribution on  $\{1, \dots, r\}$  and  $\mathcal{E}_1, \dots, \mathcal{E}_r$  are experiments with the same parameter set, then (see [30]), the  $\Pi$ -mixture,  $\sum_1^r \Pi_i \mathcal{E}_i$  of  $\mathcal{E}_1, \dots, \mathcal{E}_r$  is the experiment obtained by first observing a random index  $I$  with distribution  $\Pi$ , and then carrying out the experiment  $\mathcal{E}_I$ .

**NOTATION.** Before proceeding, we make some remarks on our use of the symbols  $\neq, \wedge, \vee, \|\cdot\|, \Lambda$  and  $\mathcal{L}$ .  $\#$  is short for "the number of elements in". The notations  $\wedge$  and  $\vee$  are used on several occasions for inf and sup. If, in particular,  $\mu_t; t \in T$  are measures then  $\wedge_t \mu_t$  and  $\vee_t \mu_t$  are, respectively, notations for  $\inf_t \mu_t$  and  $\sup_t \mu_t$  for the family  $\{\mu_t; t \in T\}$  with respect to the setwise ordering of measures. If  $\mu$  is a measure, then  $\|\mu\|$  denotes the total variation of  $\mu$ : i.e.,  $\|\mu\| = \sup\{\int f d\mu; -1 \leq f \leq 1\}$ . We will reserve the letter  $\Lambda$  as the notation for the set of all prior distributions with finite support. The letter  $\mathcal{L}$  will be used as notation for "distribution of". Thus, for example,  $\mathcal{L}_P(f)$  is the distribution of  $f$  under  $P$ .

If  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  and  $\Theta$  is finite then the distribution on the set of prior distributions on  $\Theta$  induced from  $\sum_\theta P_\theta$  by the map  $(dP_\theta/d\sum P_\theta; \theta \in \Theta)$  is called the standard measure of  $\mathcal{E}$ . This measure characterizes  $\mathcal{E}$  up to equivalence. We refer the reader to Le Cam [17] for further information on these measures.

Important functionals of experiments may be defined as follows. Let  $\phi$  be a homogenous and measurable function on  $[0, \infty[^\ominus$  and suppose  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  is dominated by the  $\sigma$ -finite measure  $\sigma$ . Then we may put  $\phi(\mathcal{E}) = \int \phi(dP_\theta; \theta \in \Theta) = \int \phi(dP_\theta/d\sigma; \theta \in \Theta) d\sigma$ .

If  $P$  and  $Q$  are probability measures then  $\int |dP - dQ|, \int (\sqrt{dP} - \sqrt{dQ})^2$  and  $\int \sqrt{dP dQ}$  are respectively: the statistical distance, the squared Hellinger distance,  $D^2(P, Q)$  and the affinity  $\gamma(P, Q)$  between them.

The *Hellinger transform* of  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  may be defined as the map  $H_\mathcal{E}$  which associates with each prior distribution  $t$  having finite support the number  $H_\mathcal{E}(t) =$

$\int \Pi_{\theta} dP_{\theta}^{\otimes n}$ . Thus  $\gamma(P, Q) = H_{\mathcal{E}}(1/2, 1/2)$ , where  $\mathcal{E} = (P, Q)$ . The Hellinger transform converts products of experiments into products of functions i.e.:

$$(2.1) \quad H_{\mathcal{E}_1 \times \dots \times \mathcal{E}_n} = H_{\mathcal{E}_1} \cdot H_{\mathcal{E}_2, \dots}, H_{\mathcal{E}_n}.$$

Let  $\mathcal{E} = ((\mathcal{X}, \mathcal{A}), (P_{\theta}, \theta \in \Theta))$  and  $\mathcal{F} = ((\mathcal{Y}, \mathcal{B}), (Q_{\theta}; \theta \in \Theta))$  be two experiments with the same parameter set  $\Theta$  and let  $\theta \rightsquigarrow \varepsilon_{\theta}$  be a nonnegative function on  $\Theta$ . We shall say, following Le Cam [14], that  $\mathcal{E}$  is  $\varepsilon$ -deficient relative to  $\mathcal{F}$  if to each finite decision space  $D$ , every family  $W_{\theta}; \theta \in \Theta$  of loss functions on  $D$  and every risk function  $r$  obtainable in  $\mathcal{F}$ , there is a risk function  $r'$  obtainable in  $\mathcal{E}$  so that:

$$(2.2) \quad r'(\theta) \leq r(\theta) + \varepsilon_{\theta} \sup_{d \in D} |W_{\theta}(d)|; \quad \theta \in \Theta$$

Restricting attention to decision spaces  $D$  where  $\#D \leq k$  we obtain the definition of  $\varepsilon$ -deficiency for  $k$ -decision problems; see [28]. If decision rules are defined as in Le Cam [14] then  $\varepsilon$ -deficiency (for  $k$ -decision problems) for all finite subsets of  $\Theta$  implies and is implied by  $\varepsilon$ -deficiency. If  $\mathcal{E}$  is 0-deficient relative to  $\mathcal{F}$  (for  $k$ -decision problems), then we shall say that  $\mathcal{E}$  is more informative than  $\mathcal{F}$  (for  $k$ -decision problems) and write this  $\mathcal{E} \geq \mathcal{F}$  ( $\mathcal{E} \geq_k \mathcal{F}$ ).

If  $\mathcal{E} \geq \mathcal{F}$  and  $\mathcal{F} \geq \mathcal{E}$  then we shall say that  $\mathcal{E}$  and  $\mathcal{F}$  are equivalent and write this  $\mathcal{E} \sim \mathcal{F}$ . The deficiency of the experiment  $\mathcal{E}$  with respect to the experiment  $\mathcal{F}$  is the greatest lower bound of all constants  $\varepsilon \geq 0$  such that  $\mathcal{E}$  is  $\varepsilon$ -deficient with respect to  $\mathcal{F}$ . This number will be denoted by  $\delta(\mathcal{E}, \mathcal{F})$ . The deficiency is not symmetric and, henceforth, not a proper distance. A distance  $\Delta$  for experiments is obtained by putting

$$(2.3) \quad \Delta(\mathcal{E}, \mathcal{F}) = \max\{\delta(\mathcal{E}, \mathcal{F}), \delta(\mathcal{F}, \mathcal{E})\};$$

see Le Cam [14]. Similarly we may define deficiencies  $\delta_k$  and distances  $\Delta_k$  based on  $k$ -decision problems.

According to Le Cam's randomization criterion, theorem 3 in [14],  $\mathcal{E} = (\mathcal{X}, \mathcal{A}, P_{\theta}; \theta \in \Theta)$  is  $\varepsilon$ -deficient with respect to  $\mathcal{F} = (\mathcal{Y}, \mathcal{B}, Q_{\theta}; \theta \in \Theta)$  if and only if there is a transition  $M$  from the band generated by  $(P_{\theta}; \theta \in \Theta)$  to the band generated by  $(Q_{\theta}; \theta \in \Theta)$  so that

$$(2.4) \quad \|P_{\theta}M - Q_{\theta}\| \leq \varepsilon_{\theta}; \theta \in \Theta.$$

If  $\mathcal{E}$  is dominated, or more generally coherent, then  $M$  may be represented as a conditional probability of  $\mathcal{B}$  given  $\mathcal{A}$  which may be regularized to a proper Markov kernel from  $(\mathcal{X}, \mathcal{A})$  to  $(\mathcal{Y}, \mathcal{B})$  when  $(\mathcal{Y}, \mathcal{B})$  is Euclidean.

If  $\Theta$  is finite then, using the notations in [23],  $\Gamma$  denotes the class of sub linear functions  $\gamma$  on  $R^{\theta}$  such that

$$\gamma(e^{\theta}) = \gamma(-e^{\theta}); \theta \in \Theta \quad \text{and} \quad \sum_{\theta} \gamma(e^{\theta}) = 1.$$

Here  $e^{\theta}$ , for each  $\theta$ , is the  $\theta$ th unit vector in  $R^{\Theta}$ , i.e.,  $e^{\theta} = (0, \dots, 1, \dots, 0)$ . The subclass of  $\Gamma$  consisting of those functions in  $\Gamma$  which are maximums of  $k$ -linear functionals will be denoted by  $\Gamma_k$ . A function  $\gamma$  will be called superlinear if  $-\gamma$  is sublinear. The deficiencies are then (the sublinear function criterion) given by

$$(2.5) \quad \delta_{(k)}(\mathcal{E}, \mathcal{F}) = \sup_{\gamma \in \Gamma_{(k)}} [\gamma(\mathcal{F}) - \gamma(\mathcal{E})].$$

Let us apply these results to deficiencies  $\delta_i(\mathcal{E})$  and  $\delta_a(\mathcal{E})$  for an experiment  $\mathcal{E} = (\mathcal{X}, \mathcal{A}, P_{\theta}; \theta \in \Theta)$ .

The randomization criterion yields directly that.

$$(2.6) \quad \delta_i(\mathcal{E}) = \min_Q \sup_{\theta} \|P_{\theta} - Q\| \quad \text{and} \quad \delta_a(\mathcal{E}) = 2 \min_M \sup_{\theta} P_{\theta}M(\{\theta\}^c),$$

where  $Q$  runs through all probability distributions on  $\mathcal{A}$ , while  $M$  runs through all transitions from the band generated by the  $P_{\theta}$ 's to the band of discrete finite measures on  $\Theta$ .

It follows that  $\delta_i(\mathcal{E})$  is the minimax risk in the estimation problem where no observations are available and loss is measured by statistical distance. Furthermore

$$(2.7) \quad \frac{1}{2} \sup_{\theta_1, \theta_2} \|P_{\theta_1} - P_{\theta_2}\| \leq \delta_i(\mathcal{E}) \leq \sup_{\theta_1, \theta_2} \|P_{\theta_1} - P_{\theta_2}\|,$$

where the left-hand side is  $\delta_2(\mathcal{M}_i, \mathcal{E})$ , by [23].

Similarly  $\frac{1}{2}\delta_a(\mathcal{E})$  is the minimax risk in the problem of estimating  $\theta$  on the basis of  $\mathcal{E}$  when the loss is 1 or 0 as the estimator hits or not. If  $\mathcal{E}$  is dominated then  $(P_\theta M; \theta \in \Theta)$  is also dominated, so that  $\inf_\theta P_\theta M(\{\theta\}) = 0$  when  $\Theta$  is uncountable. Hence  $\delta_a(\mathcal{E}) = 2$  in this case. More generally  $\mathcal{E}$  is  $\epsilon$ -deficient with respect to  $\mathcal{M}_a$  if and only if  $P_\theta M(\{\theta\}) \geq 1 - \epsilon_\theta/2$ ;  $\theta \in \Theta$  for some transition  $M$ . If  $\Theta$  is finite and each  $P_\theta$  is atomless then, by Dvoretzky et al's extension [9] of Lyapunov's theorem [20], this is equivalent to the existency of a measurable partition  $A_\theta$ ;  $\theta \in \Theta$  of  $\chi$  such that  $P_\theta(A_\theta) \geq 1 - \epsilon_\theta/2$ .

If  $\lambda$  is a prior distribution on  $\Theta$  with countable support then the minimum Bayes risk in the last estimation problem will be denoted by  $b(\lambda | \mathcal{E})$ . It is easily seen that

$$(2.8) \quad b(\lambda | \mathcal{E}) = 1 - \|\vee_\theta \lambda_\theta P_\theta\|$$

and the map  $\lambda \sim \rightarrow b(\lambda | \mathcal{E})$  on  $\Lambda$  defines, by Morse and Sacksteder [21],  $\mathcal{E}$  up to equivalence. If  $\Theta$  is countable then each decision rule  $\delta$  may be represented as a random distribution which to each  $x$  assigns the distribution  $\delta.(x)$  on  $\Theta$ . Suppose in addition that each  $P_\theta$  has density  $f_\theta$  with respect to the measure  $\sigma$ . Then  $\delta$  achieves minimum Bayes risk if and only if, for  $\sigma$  almost all  $x$ ,  $\delta.(x)$  is supported by  $\{\theta; \lambda_\theta f_\theta(x) = \max_{\theta'} \lambda_{\theta'} f_{\theta'}(x)\}$ . In particular any maximum likelihood estimator of  $\theta$  achieves minimum Bayes risk when  $\theta$  is finite and  $\lambda$  is uniform. It follows from straight-forward minimax theory that

$$(2.9) \quad \frac{1}{2}\delta_a(\mathcal{E}) = \sup_{\lambda \in \Lambda} b(\lambda | \mathcal{E}) = 1 - \inf_{\lambda \in \Lambda} \|\vee_\theta \lambda_\theta P_\theta\|.$$

The prior distribution  $\lambda$  (with countable support) will be called least favorable if it is least favorable in this estimation problem, i.e., if  $b(\lambda | \mathcal{E}) = \frac{1}{2}\delta_a(\mathcal{E})$ .

It follows from [23] that  $\mathcal{E}$  is  $\epsilon$ -deficient with respect to  $\mathcal{M}_a$  for testing problems if and only if to each subset  $\Theta_0$  of  $\Theta$  there corresponds a power function  $\pi$  in  $\mathcal{E}$  so that  $\pi(\theta) \leq \epsilon_\theta/2$  or  $\geq 1 - \epsilon_\theta/2$  as  $\theta \in \Theta_0$  or  $\theta \notin \Theta_0$ . The sublinear function criterion (2.5) shows that

$$(2.10) \quad \delta_2(\mathcal{E}, \mathcal{M}_a) = \sup(1 - \|\sum_\theta a_\theta P_\theta\|)$$

where  $a$  runs through all finite measures  $a$  on  $\theta$  with finite support and total variation =  $\sum_\theta |a_\theta| = 1$ .

By [23],  $\delta_i(\mathcal{M}_a) = \delta_a(\mathcal{M}_i) = 2(1 - m^{-1})$  where  $m = \#\Theta$ . Hence, by the triangular inequality for deficiencies

$$(2.11) \quad 2 - 2/m \leq \delta_i(\mathcal{E}) + \delta_a(\mathcal{E}).$$

One might expect that  $\delta_a(\mathcal{E})$  is small (large) when  $\delta_i(\mathcal{E})$  is large (small) and conversely. If  $\Theta$  is finite, then this may be made precise by inequalities; see [29]. If  $\Theta$  is infinite, however, then it may easily happen that  $\delta_i(\mathcal{E}^n) \rightarrow 2$  while  $\delta_a(\mathcal{E}^n) \geq 1$  for all  $n$ .

**3. Replicated dichotomies.** We shall in this section assume that  $\Theta = \{1, 2\}$  i.e., that our experiments are dichotomies. In this case deficiencies may, [23], be expressed in terms of testing problems only. We shall need:

**THEOREM 3.1.** *Suppose  $\mathcal{E} = (P_1, P_2)$  and  $\mathcal{F} = (Q_1, Q_2)$ . Denote by  $\beta(\alpha | \mathcal{E})$  and  $\beta(\alpha | \mathcal{F})$  the powers of the most powerful level  $\alpha$  tests for  $\theta = 1$  against  $\theta = 2$  in  $\mathcal{E}$  and  $\mathcal{F}$ . Also denote by  $b(\lambda | \mathcal{E})$  and  $b(\lambda | \mathcal{F})$  the minimum Bayes risks for the same problem for 0 - 1 loss and prior  $(\lambda_1, \lambda_2)$  in  $\mathcal{E}$  and  $\mathcal{F}$ .*

*Then each of the following conditions are equivalent:*

$$(i) \quad \beta\left(\alpha + \frac{\varepsilon_1}{2} \mid \mathcal{E}\right) + \frac{\varepsilon_2}{2} \cong \beta(\alpha \mid \mathcal{F}) \quad ; \alpha \cong 0$$

$$(ii) \quad [\lambda_1 \varepsilon_1 + \lambda_2 \varepsilon_2] / 2 \cong b(\lambda \mid \mathcal{E}) - b(\lambda, \mathcal{F}) \quad ; 0 \cong \lambda \cong 1$$

$$(iii) \quad \int \gamma(dP_2/dP_1) \cong \int \gamma(dQ_2/dQ_1) dQ_1 + \frac{\varepsilon_1}{2} [\gamma(\infty) - \gamma(0)] + \frac{\varepsilon_2}{2} \gamma'(0)$$

for each concave function  $\gamma$  on  $[0, \infty[$  such that  $\gamma(x)/x \rightarrow 0$  as  $x \rightarrow \infty$ .

$$(iv) \quad \frac{\varepsilon_1}{2} [\phi'(1) - (\phi(1) - \phi(0))] + \frac{\varepsilon_2}{2} [(\phi(1) - \phi(0)) - \phi'(0)] \\ \cong \int \phi(dQ_2/d(Q_1 + Q_2))d(Q_1 + Q_2) - \int \phi(dP_2/d(P_1 + P_2))d(P_1 + P_2)$$

for any convex function  $\phi$  on  $[0, 1]$ .

(v)  $\mathcal{E}$  is  $(\varepsilon_1, \varepsilon_2)$  deficient with respect to  $\mathcal{F}$ .

**REMARK.** If  $\mu$  and  $\nu$  are finite measures then  $d\mu/d\nu$  are the Radon-Nikodym derivative of the  $\nu$  continuous part of  $\mu$  with respect to  $\nu$ .

**PROOF.** The equivalence of (i), (ii) and (v) follows from [23] while the implications (iv)  $\Rightarrow$  (iii)  $\Rightarrow$  (ii) follows by applying (iv) and (iii) to, respectively, the functions  $\phi(x) = \gamma\left(\frac{x}{1-x}\right)(x-1)$  and  $\gamma(x) = \min\{1-\lambda, \lambda x\}$ . It remains therefore to show that (ii) implies (iv). Suppose (ii) holds and let  $\phi$  be convex on  $[0, 1]$ . Define convex functions  $\phi_n; n=1, 2, \dots$  by requiring that:  $\phi_n(k/2^n) = \phi(k/2^n); k=0, 1, \dots, 2^n$  and that  $\phi_n$  is linear on each interval  $\left[\frac{i-1}{2^n}, \frac{i}{2^n}\right]; i=1, \dots, 2^n$ . Then  $\phi_n \downarrow \phi, \phi'_n(0) \downarrow \phi'(0)$  and  $\phi'_n(1) \uparrow \phi'(1)$ . It follows, without loss of generality, that we may assume that the graph of  $\phi$  consists of a finite number of line segments. Hence we may write  $\phi(x) = \max\{a_i + b_i x; i=1, \dots, k\}$  where  $a_1 > a_2 > \dots > a_k$  and  $b_1 < b_2 < \dots < b_k$ . It follows that:  $\phi(x) = a_1 + b_1 x + \sum_{i=1}^k (a_{i+1} + b_{i+1} x - a_i - b_i x)^+$ . The inequality is trivial when  $\phi$  is linear. It suffices therefore to prove it for functions  $\phi$  of the form

$$\phi(x) = A \vee Bx = Bx + A(1-x) - \min\{A(1-x), (B-A)x\}$$

where, since the case  $A \cong B$  is trivial,  $0 < A < B$ . Putting  $\lambda = 1 - A/B$  we see then that the inequality reduces to (ii).  $\square$

A prior distribution  $\lambda$  may in this case be identified with the probability it assigns to  $\{2\}$ . We shall therefore in this Section write  $b(\lambda \mid \mathcal{E})$  instead of  $b(1-\lambda, \lambda \mid \mathcal{E})$  and  $H(t \mid \mathcal{E})$  instead of  $H(1-t, t \mid \mathcal{E})$  when  $\lambda, t \in [0, 1]$ . The Bayes risk function  $\lambda \rightsquigarrow b(\lambda \mid \mathcal{E})$  may then be written

$$b(\lambda \mid \mathcal{E}) = 1 - \|(1-\lambda)P_1 \vee \lambda P_2\| = \|(1-\lambda)P_1 \wedge \lambda P_2\|.$$

Hence

$$0 \cong b(\lambda \mid \mathcal{E}) \cong (1-\lambda) \wedge \lambda; \quad \lambda \in [0, 1].$$

Note that the right-hand side is  $b(\lambda \mid \mathcal{M}_i)$  while the left-hand side is  $b(\lambda \mid \mathcal{M}_a)$ . Conversely, any concave function  $f$  on  $[0, 1]$  such that  $0 \cong f(\lambda) \cong (1-\lambda) \wedge \lambda$  for  $\lambda \in [0, 1]$  is of the form  $f(\lambda) = b(\lambda \mid \mathcal{E})$ , where  $\mathcal{E}$  is, up to equivalence, determined by  $f$ . A particularly interesting aspect of this representation is the relation

$$b(\lambda \mid \sup_t \mathcal{E}_t) = \inf_t b(\lambda \mid \mathcal{E}_t),$$

which is valid for any family  $\{\mathcal{E}_t; t \in T\}$  of dichotomies. Deficiencies are easily described

in terms of these functions. By Theorem 3.1

$$\delta(\mathcal{E}, \mathcal{F}) = 2 \sup_{0 \leq \lambda \leq 1} [b(\lambda | \mathcal{E}) - b(\lambda | \mathcal{F})]$$

for any pair  $(\mathcal{E}, \mathcal{F})$  of dichotomies. In particular  $\delta_i(\mathcal{E}) = 1 - 2b(\frac{1}{2}, \mathcal{E})$ ; see [23].

There is a simple connection between the Hellinger transform and the Bayes risk for dichotomies as follows.

**THEOREM 3.2.** [28]

$$H(t | \mathcal{E}) / (1 - t)t = \int_0^1 b(\lambda | \mathcal{E})(1 - \lambda)^{t-2}\lambda^{-t-1}d\lambda = \int_0^\infty x^{-t-1} \|xP_1 \wedge P_2\| dx.$$

As an application consider the problem of finding bounds for  $H(t | \mathcal{E})$  in terms of  $b(\lambda | \mathcal{E})$ .

**THEOREM 3.3.** For any pair  $(\lambda, t)$  of prior distributions.

$$b(\lambda | \mathcal{E}) \leq (1 - \lambda)^{1-t}\lambda^t H(t | \mathcal{E})$$

and

$$H(t | \mathcal{E}) \leq [1 - b(\lambda | \mathcal{E}) + \lambda b'(\lambda | \mathcal{E})]^{1-t} [b(\lambda | \mathcal{E}) + (1 - \lambda)b'(\lambda | \mathcal{E})]^t + [b(\lambda | \mathcal{E}) - \lambda b'(\lambda | \mathcal{E})]^{1-t} [1 - b(\lambda | \mathcal{E}) - (1 - \lambda)b'(\lambda | \mathcal{E})]^t,$$

where  $b'(\lambda | \mathcal{E})$  is any number between the left and right derivative of  $\lambda \rightsquigarrow b(\lambda | \mathcal{E})$  at  $\lambda$ . Both inequalities are, for given values of  $b(\lambda | \mathcal{E})$ , sharp.

**PROOF.** Put  $\lambda = \lambda_0$ ,  $b_0 = b(\lambda_0 | \mathcal{E})$ ,  $b'_0 = b'(\lambda_0 | \mathcal{E})$  and suppose  $0 < \lambda_0 < 1$ . Then  $b(\cdot | \mathcal{C}) \leq b(\cdot | \mathcal{E}) \leq b(\cdot | \mathcal{D})$  where  $b(\lambda | \mathcal{C}) = b_0[(\lambda/\lambda_0) \wedge (1 - \lambda)/(1 - \lambda_0)]$  and  $b(\lambda | \mathcal{D}) = [b_0 + (\lambda - \lambda_0)b'_0] \wedge (1 - \lambda) \wedge \lambda$ . Hence  $\mathcal{C} \cong \mathcal{E} \cong \mathcal{D}$  so that  $H(t | \mathcal{C}) \leq H(t | \mathcal{E}) \leq H(t | \mathcal{D})$ . These are the desired inequalities.  $\square$

**REMARK.** The first inequality implies that

$$b(\lambda | \mathcal{E}) \leq ((1 - \lambda) \wedge \lambda) \inf_t H(t | \mathcal{E}),$$

yielding the right of the inequalities:

$$[\gamma(\mathcal{E})/2]^2 \leq \delta_a(\mathcal{E})/2 \leq \inf_t H(t | \mathcal{E}),$$

while the left inequality follows by choosing  $\lambda$  least favorable and by putting  $t = \frac{1}{2}$ .

Inserting  $\lambda = t = \frac{1}{2}$  we find

$$\|P_1 \wedge P_2\| \leq \gamma(P_1, P_2) \leq \sqrt{\|P_1 \wedge P_2\| \|P_1 \vee P_2\|},$$

or equivalently

$$D^2(P_1, P_2) \leq \|P_1 - P_2\| \leq 2\sqrt{1 + (P_1, P_2)} D(P_1, P_2).$$

This proves the known fact, see Le Cam [15], that the Hellinger distance is equivalent to the statistical distance.

Using the fact that

$$\lambda_1 f_1 g_1 \wedge \lambda_2 f_2 g_2 \cong (\lambda_1 f_1 \wedge \lambda_2 f_2)(g_1 \wedge g_2) \cong (\lambda_1 f_1 \wedge \lambda_2 f_2)(\mu_1 g_1 \wedge \mu_2 g_2)$$

whenever  $\lambda_1, \lambda_2, \mu_1, \mu_2, f_1, f_2, g_1$  and  $g_2$  are nonnegative numbers we get:

**PROPOSITION 3.4.** *If  $\mathcal{E}$  and  $\mathcal{F}$  are dichotomies, then for any pair  $(\lambda, \mu)$  of prior distributions*

$$b(\lambda | \mathcal{E} \times \mathcal{F}) \geq 2b(\lambda | \mathcal{E})b(\frac{1}{2} | \mathcal{F}) \geq b(\lambda | \mathcal{E})b(\mu | \mathcal{F}).$$

*In particular*

$$\delta_a(\mathcal{E} \times \mathcal{F})/2 \geq [\delta_a(\mathcal{E})/2]2b(\frac{1}{2} | \mathcal{F}) \geq [\delta_a(\mathcal{E})/2][\delta_a(\mathcal{F})/2].$$

By proposition 3.4,  $b(\lambda | \mathcal{E}^{m+n}) \geq b(\lambda | \mathcal{E}^m)b(\lambda | \mathcal{E}^n)$  and  $\delta_a(\mathcal{E}^{m+n})/2 \geq [\delta_a(\mathcal{E}^m)/2][\delta_a(\mathcal{E}^n)/2]$ . It follows that  $\delta_a(\mathcal{E}^n)^{1/n}, n = 1, 2, \dots$  and  $b(\lambda | \mathcal{E}^n)^{1/n}, n = 1, 2, \dots$  converge, as  $n \rightarrow \infty$ , to respectively  $\sup_n [\delta_a(\mathcal{E}^n)/2]^{1/n}$  and  $\sup_n b(\lambda | \mathcal{E}^n)^{1/n}$ . Furthermore, since  $b(\lambda | \mathcal{E}^n) \leq \delta_a(\mathcal{E}^n)/2 \leq b(\lambda | \mathcal{E}^n)/[(1 - \lambda) \wedge \lambda]$  when  $\lambda \in ]0, 1]$ , these limits are the same. Thus

**THEOREM 3.5.** *There is for each dichotomy  $\mathcal{E}$ , a constant  $C(\mathcal{E})$  in  $[0, 1]$  such that, for each nondegenerate prior distribution  $\lambda$*

$$\begin{aligned} \lim_{n \rightarrow \infty} b(\lambda | \mathcal{E}^n)^{1/n} &= \lim_{n \rightarrow \infty} [\delta_a(\mathcal{E}^n)/2]^{1/n} = \sup_n b(\lambda | \mathcal{E}^n)^{1/n} \\ &= \sup_n [\delta_a(\mathcal{E}^n)/2]^{1/n} = C(\mathcal{E}). \end{aligned}$$

**REMARK.** It follows that  $C(\mathcal{E}^r) = C(\mathcal{E})^r; r = 1, 2, \dots$  and that  $C(\mathcal{E}) \geq C(\mathcal{F})$  when  $\mathcal{E} \leq \mathcal{F}$ . Also, by this Theorem and Proposition 3.4:  $C(\mathcal{E} \times \mathcal{F}) \geq C(\mathcal{E})C(\mathcal{F})$  for any pair  $(\mathcal{E}, \mathcal{F})$  of dichotomies. Using that  $1 - \delta_i(\mathcal{E}) = 2b(\frac{1}{2}, \mathcal{E})$  we get

**COROLLARY 3.6.**

$$\lim_{n \rightarrow \infty} [1 - \delta_i(\mathcal{E}^n)]^{1/n} = C(\mathcal{E}).$$

By the remark after Theorem 3.3:  $[\delta_a(\mathcal{E}^n)/2]^{1/n} \leq \inf_t H(t | \mathcal{E})$ . The fact that the limit  $C(\mathcal{E})$  of the left-hand side equals the right-hand side is a consequence of the following result of Chernoff [6].

**THEOREM 3.7.** *For any dichotomy  $\mathcal{E}$  and any nondegenerate prior distribution  $\lambda$*

$$\lim_{n \rightarrow \infty} b(\lambda | \mathcal{E}^n)^{1/n} = \inf_t H(t | \mathcal{E}).$$

Chernoff first derives the basic convergence result for large deviations,

$$\lim_{n \rightarrow \infty} P(1/n (X_1 + \dots + X_n) \geq 0)^{1/n} = \inf_{t \geq 0} Ee^{tX},$$

which is valid for any sequences  $X_1, X_2, \dots$  of independent and identically distributed variables. Then, applying this to log likelihoods, he proves Theorem 6.7.

Chernoff's result has simple and important interpretations in terms of experiments and it is somewhat unsatisfactory that several of the steps in the various proofs are not so easily interpreted in terms of experiments. We shall now give a simple and "natural" proof of this important result which is based solely on a few basic facts about statistical experiments. The proof admits variations and is of considerable interest in itself. We shall see in the next section (Theorem 7.2) how this idea may be used to extend Chernoff's result to a more general result. The main tools of the proof are compactness of  $\Delta$  convergence, [17], and a type of product homomorphism for experiments which we shall describe now.

Let  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  be an experiment dominated by the  $\sigma$ -finite measure  $\mu$ . Put  $f_\theta = dP_\theta/d\mu$  and consider a family  $\{t^{(i)}; i \in I\}$  of prior distributions with finite supports and



such that  $H(t^{(i)} | \mathcal{E}) > 0, i \in I$ . For each  $i \in I$  define a probability measure  $Q_i$  by  $dQ_i/d\mu = H(t^{(i)} | \mathcal{E})^{-1} \prod_{\theta} f_{\theta}^{t_{\theta}^{(i)}}$ . It is easily seen that the measures  $Q_i$ , do not depend on our choice of dominating measure. We shall denote the experiment  $(Q_i: i \in I)$  by  $(\mathcal{E}: t^{(i)}; i \in I)$ . If  $s$  is a prior on  $I$  with finite support, then

$$H(s | (\mathcal{E}: t^{(i)}; i \in I)) = H(\sum_i s_i t^{(i)} | \mathcal{E}) \prod_i H(t^{(i)} | \mathcal{E})^{-s_i}.$$

It follows that  $(\mathcal{E}: t^{(i)}; i \in I) \sim (\mathcal{F}: t^{(i)}; i \in I)$  when  $\mathcal{E} \sim \mathcal{F}$ . Also, by the product rule for Hellinger transforms

$$(\mathcal{E}: t^{(i)}; i \in I) \times (\mathcal{F}: t^{(i)}; i \in I) \sim (\mathcal{E} \times \mathcal{F}: t^{(i)}; i \in I)$$

in the sense that the left side is defined if and only if the right side is defined and then equivalence holds.

Consider now a dichotomy  $\mathcal{E} = (P_1, P_2) \sim \mathcal{M}_a$  and a real number  $a$ . Put  $\rho(t) = e^{-at}H(t | \mathcal{E})$  and  $\rho_0 = \inf_t \rho(t)$  and assume  $\rho_0 = \rho(t_0)$  where  $t_0 \in ]0, 1[$ . Then, when  $n$  is sufficiently large, the Hellinger transform of

$$\left( \mathcal{E}: t_0; t_0 + \frac{s}{\sqrt{n}} \right)^n \sim \left( \mathcal{E}^n: t_0; t_0 + \frac{s}{\sqrt{n}} \right)$$

is

$$\left[ \rho \left( t_0 + \frac{ts}{\sqrt{n}} \right) / \rho_0^{1-t} \rho \left( t_0 + \frac{s}{\sqrt{n}} \right) \right]^{t^n}$$

which converges, as  $n \rightarrow \infty$ , to  $e^{-s^2\tau^2 t(1-t)/2}$  where  $\tau^2 = \rho''(t_0)/\rho_0$ . This limit, as a function of  $t$ , is the Hellinger transform of the dichotomy  $\mathcal{G}_s = (Q_0, Q_s)$  where  $Q_s$ , for each  $s$ , is the Normal  $(s\tau, 1)$  distribution. Hence  $\Delta \left( \left( \mathcal{E}^n: t_0; t_0 + \frac{s}{\sqrt{n}} \right), \mathcal{G}_s \right) \rightarrow 0$  so that the laws of

likelihood ratios of  $\left( \mathcal{E}^n: t_0; t_0 + \frac{s}{\sqrt{n}} \right)$  converges weakly to the corresponding laws for  $\mathcal{G}_s$ .

Thus,

$$\mathcal{L} \left( U_n^{s/\sqrt{n}} e^{-as\sqrt{n}} \left( t_0 + \frac{s}{\sqrt{n}} \right)^{-n} \rho_0^n \mid e^{-nat_0} \rho_0^{-n} U_n^t P_1^n \right) \rightarrow \mathcal{L}(U \mid Q_0),$$

where  $U_n = dP_2^n/dP_1^n$  and  $U = dQ_s/dQ_0$ . Hence, since  $\rho \left( t_0 + \frac{s}{\sqrt{n}} \right)^{-n} \rho_0^n \rightarrow e^{-s^2\tau^2/2}$ ,

$$E_{P_1^n} \phi(U_n^{s/\sqrt{n}} e^{-as\sqrt{n}}) U_n^{t_0} \sim \rho_0^n e^{nat_0} E_{Q_0} \phi(e^{s^2\tau^2/2} U)$$

when  $\phi$  is bounded and continuous a.e. Lebesgue if  $\tau > 0$  and continuous in 1 if  $\tau = 0$ .

Choose a number  $M > 1$  and put  $\phi(x) = 1$  or  $0$  as  $x \in [M^{-1}, M]$  or not.

Then

$$\begin{aligned} \| P_1^n \wedge e^{-na} P_2^n \| &\geq \int (1 \wedge e^{-na} U_n) \phi(e^{-\sqrt{na}} U_n^{1/\sqrt{n}}) U_n^{-t_0} U_n^{t_0} dP_1^n \\ &\geq M^{-\sqrt{n}} M^{-\sqrt{n}t_0} e^{-nat_0} \int \phi(e^{-\sqrt{na}} U_n^{1/\sqrt{n}}) U_n^{t_0} dP_1^n \\ &\sim M^{-\sqrt{n}(1+t_0)} \rho_0^n E_{Q_0} \phi(e^{\tau^2/2} dQ_1/dQ_0). \end{aligned}$$

Hence, since the last expectation is  $> 0$ , we get

$$(3.1) \quad \liminf_n \| P_1^n \wedge e^{-na} P_2^n \|^{1/n} \geq \inf_t e^{-at} H(t | \mathcal{E}).$$

If  $\mathcal{E} \sim \mathcal{M}_a$  then  $\rho_0 = 0$  so that (3.1) holds trivially. In general the function  $\rho$  may be forced

to obtain its infimum in ]0, 1[ by increasing the amount of information in  $\mathcal{E}$  by a negligible amount. To be more precise let  $N$  be large and put  $\mathcal{G} = \left(1 - \frac{1}{N}\right) \mathcal{M}_i + \frac{1}{N} \mathcal{F} = (\bar{P}_1, \bar{P}_2)$  where  $-H'(0 + | \mathcal{F}) = H'(1 - | \mathcal{F}) = \infty$ . Then  $t \rightsquigarrow H(t | \mathcal{E} \times \mathcal{G})e^{-at}$  obtains its infimum in ]0, 1[. Hence

$$\begin{aligned} \liminf_n \| P_1^n \wedge e^{-na} P_2^n \|^{1/n} &\geq \liminf_n \| (P_1 \times Q_1)^n \wedge e^{-n} (P_2 \times Q_2)^n \| \\ &\geq \inf_t [e^{-at} H(t | \mathcal{E}) H(t | \mathcal{G})] \geq \left(1 - \frac{1}{N}\right) \rho_0 \rightarrow \rho_0 \end{aligned}$$

as  $N \rightarrow \infty$ . We used here that  $H(t | \mathcal{G}) = \left(1 - \frac{1}{N}\right) + \frac{1}{N} H(t | \mathcal{F}) \geq 1 - \frac{1}{N}$ .

It follows that (3.1) holds for any dichotomy  $(P_1, P_2)$ .  $a = 0$  yields  $\liminf_n b(\lambda | \mathcal{E}^n)^{1/n} \geq \inf_t H(t | \mathcal{E})$  while  $b(\lambda | \mathcal{E}^n) \leq H(t | \mathcal{E})^n$  by Theorem 3.3. Hence  $b(\lambda | \mathcal{E}^n)^{1/n} \rightarrow \inf_t H(t | \mathcal{E})$  and this is Chernoff's theorem. By Theorems 3.5 and 3.7

**THEOREM 3.8.**

$$C(\mathcal{E}) = \inf_t H(t | \mathcal{E}).$$

**REMARK.** By the remarks after Theorem 3.3,  $b(\lambda | \mathcal{E}) \leq [(1 - \lambda) \vee \lambda] C(\mathcal{E})$  and  $\gamma[(\mathcal{E})/2]^2 \leq \delta_a(\mathcal{E})/2 \leq C(\mathcal{E})$ .

We shall now consider a few extensions of Theorem 3.5. Let us first consider the asymptotic behavior of minimum Bayes risk in other decision problems. It is known, see for example [8] or [23], that the minimum Bayes risk may often be expressed as functionals  $\psi(\mathcal{E}) = \int \psi(dP_1, dP_2)$  where the function  $\psi$  is super linear (i.e.,  $\psi(x + y) \geq \psi(x) + \psi(y)$  and  $\psi(tx) = t\psi(x)$  when  $t \geq 0$ ) or  $R^2$ . The function  $\psi$  is determined by the loss function. It follows, since the standard measure of  $\mathcal{E}^n$  converges weakly to the standard measure of  $\mathcal{M}_a$  when  $\mathcal{E} \not\sim \mathcal{M}_i$ , that  $\psi(\mathcal{E}^n) \rightarrow \psi(\mathcal{M}_a)$  as  $n \rightarrow \infty$  provided  $\mathcal{E} \not\sim \mathcal{M}_i$ . By Theorem 2 in [23]:

$$0 \leq \psi(\mathcal{E}) - \psi(\mathcal{M}_a) \leq \frac{\delta_a(\mathcal{E})}{2} [\psi(1, 0) + \psi(-1, 0) + \psi(0, 1) + \psi(0, -1)].$$

It follows, by replacing  $\mathcal{E}$  with  $\mathcal{E}^n$  and applying Theorem 3.5 that  $\limsup_n [\psi(\mathcal{E}^n) - \psi(\mathcal{M}_a)]^{1/n} \leq C(\mathcal{E})$ .

Suppose  $\psi$  is not affine on  $[0, \infty[^2$ . Put  $\phi(x) = \psi(1 - x, x)$ ,  $x \in [0, 1]$ . Then  $\phi$  is concave on  $[0, 1]$  and for some  $x_0 \in ]0, 1[$ :  $\phi(x_0) > (1 - x_0)\phi(0) + x_0\phi(1)$ . Let  $\chi$  be the function on  $[0, 1]$  which is linear on the intervals  $[0, x_0]$  and  $[x_0, 1]$  and which satisfies  $\chi(0) = \phi(0)$ ,  $\chi(x_0) = \phi(x_0)$ ,  $\chi(1) = \phi(1)$ . Then  $\psi(\mathcal{E}) - \psi(\mathcal{M}_a) = \tilde{\psi}(\mathcal{E}) - \tilde{\psi}(\mathcal{M}_a)$  where  $\tilde{\psi}(x) = \psi(x) - x_1\psi(1, 0) - x_2\psi(0, 1)$ . Thus we may as well assume that  $\phi(0) = \phi(1) = 0$  and then  $\psi(\mathcal{E}) - \psi(\mathcal{M}_a) = \psi(\mathcal{E}) \geq \int \chi \left( \frac{dP_2}{d(P_1 + P_2)} \right) d(P_1 + P_2) = kb(\lambda, \mathcal{E})$  where  $k = \phi(x_0) \left( \frac{1}{x_0} + \frac{1}{1 - x_0} \right)$  and  $\lambda = \frac{1}{x_0} \left( \frac{1}{x_0} + \frac{1}{1 - x_0} \right)^{-1}$ . Using Theorem 3.5 once more, we find the following.

**THEOREM 3.9.** *If  $\psi$  is super linear on  $R^2$  and not linear on  $[0, \infty[^2$  then*

$$\lim_{n \rightarrow \infty} [\psi(\mathcal{E}^n) - \psi(\mathcal{M}_a)]^{1/n} = C(\mathcal{E}).$$

**REMARK.** If  $\psi(x) = \wedge_t \sum_{\theta=1}^2 \lambda_\theta L_\theta(t) x_\theta$  where  $L$  is the loss function, then  $\psi(\mathcal{E})$  is the minimum Bayes risk for the loss function  $L$ . If  $\wedge_t$  is attained, then the exceptional case is the situation where, for some  $t_0$ ,  $L_\theta(t_0) \leq L_\theta(t)$  for  $\theta = 1, 2$ , and all  $t$ . In that case no observations are needed and the decision rule  $x \rightsquigarrow t_0$  is "uniformly" optimal.

In spite of the last remark there are interesting measures of information based on

deficiencies whose exponential rates of convergence differs from  $C(\mathcal{E})$ .

**EXAMPLE 3.10.** Fix a number  $\alpha \in [0, 1]$  and consider the smallest number  $2\epsilon_2$  such that the dichotomy  $(P_1, P_2)$  is  $(2\alpha, 2\epsilon_2)$  deficient with respect to  $\mathcal{M}_\alpha$ . By Theorem 2.3 this is just  $1 - \beta(\alpha | P_1, P_2)$  where  $\beta(\alpha | P_1, P_2)$  is the power of the most powerful level  $\alpha$  test for testing “ $P_1$ ” against “ $P_2$ ”. It is shown in an unpublished paper of Stein (see [7] or [1] for a proof) that

$$[1 - \beta(\alpha | P_1^n, P_2^n)]^{1/n} \rightarrow e^{E_n \log \frac{dP_2}{dP_1}}$$

We shall for the remaining part of this section assume that  $\mathcal{L}_{P_1}(\log(dP_2/dP_1))$  is nonlattice and that  $\inf H(t | \mathcal{E})$  is obtained at  $t_0 \in ]0, 1[$ . Put  $\tau^2 = C(\mathcal{E})^{-1}H''(t_0 | \mathcal{E})$  and  $A_n = C(\mathcal{E})^{-n}(2\pi\tau^2n)^{1/2}$ . Then the expansion of Efron and Truax [11] may be written

$$b(\lambda | \mathcal{E}^n) \sim A_n^{-1}(1 - \lambda)^{1-t_0} \lambda^{t_0}/(1 - t_0)t_0.$$

The functions  $\lambda \rightsquigarrow b(\lambda | \mathcal{E}^n)$  and the function  $\lambda \rightsquigarrow (1 - \lambda)^{1-t_0} \lambda^{t_0}$  are all concave on  $[0, 1]$ . It follows that the convergence is uniform in  $\lambda$ . Maximizing with respect to  $\lambda$  we find that  $\delta_\alpha(\mathcal{E}^n)/2 \sim A_n^{-1}(1 - t_0)^{-t_0}t_0^{t_0-1}$  i.e.,  $\delta_\alpha(\mathcal{E}^n)/2 = b(t_0 | \mathcal{E}^n)(1 + o(1))$ . Let  $\mu_n$  be least favorable in  $\mathcal{E}^n$ , i.e.,  $b(\mu_n | \mathcal{E}^n) = \delta_\alpha(\mathcal{E}^n)/2$ . Then

$$A_n b(\mu_n, \mathcal{E}^n) - [(1 - \mu_n)^{1-t_0} \mu_n^{t_0}/(1 - t_0)t_0] \rightarrow 0.$$

Hence  $(1 - \mu_n)^{1-t_0} \mu_n^{t_0} \rightarrow (1 - t_0)^{1-t_0} t_0^{t_0}$  so that  $\mu_n \rightarrow t_0$ . By a slight extension of this argument we find that  $\delta_\alpha(\mathcal{E}^n)/2 = b(\mu_n | \mathcal{E})(1 + o(1))$  as  $n \rightarrow \infty$ , if and only if  $\mu_n \rightarrow t_0$ . This proves

**THEOREM 3.11.**

- (i)  $\delta_\alpha(\mathcal{E}^n)/2 = (1 - t_0)^{-t_0} t_0^{t_0-1} \frac{1}{\sqrt{2\pi\tau^2}} \frac{1}{\sqrt{n}} (1 + o(1))$  as  $n \rightarrow \infty$ .
- (ii)  $t_0$  is asymptotically least favorable in the sense that  $\delta_\alpha(\mathcal{E}^n)/2 = b(t_0 | \mathcal{E}^n)(1 + o(1))$  as  $n \rightarrow \infty$ .
- (iii) More generally  $\delta_\alpha(\mathcal{E}^n)/2 = b(\mu_n | \mathcal{E}^n)(1 + o(1))$  as  $n \rightarrow \infty$  if and only if  $\lim_n \mu_n = t_0$ .

The prior  $t_0$  which minimizes  $t \rightsquigarrow H(t | \mathcal{E})$  is, by Theorem 3.11, asymptotically least favorable.

Let  $U$  be the distribution on  $]0, \infty[$  with density  $x \rightsquigarrow x^{-t_0}$  with respect to Lebesgue measure. Note that the expansion of Efron and Truax may be written,

$$\lim_{n \rightarrow \infty} \int [(1 - \lambda) \wedge \lambda x] A_n K_n(dx) = \int [(1 - \lambda) \wedge \lambda x] U(dx)$$

where  $K_n = \mathcal{L}_{P_1^n}(dP_2^n/dP_1^n)$ . It follows that  $\lim_{n \rightarrow \infty} \int \phi(x) A_n K_n(dx) = \int \phi(x) U(dx)$  for any function  $\phi$  on  $[0, \infty[$  which is a linear combination of functions  $x \rightsquigarrow (1 - \lambda) \wedge \lambda x$ ;  $\lambda \in [0, 1]$ . It is not difficult to see that a function  $\phi$  is a linear combination of functions  $x \rightsquigarrow (1 - \lambda) \wedge \lambda x$ ;  $\lambda \in [0, 1]$ , if and only if  $\phi$  is polygonal,  $\phi(0) = 0$  and  $\phi(x) = \lim_{x \rightarrow \infty} \phi(x)$  when  $x$  is sufficiently large. Hence, by the theory of weak convergence of measures,

**THEOREM 3.12.**

$$\int \phi(dP_2^n/dP_1^n) dP_1^n = \left[ \int \phi(x) U_{t_0}(dx) \right] \frac{1}{\sqrt{2\pi\tau^2}} \frac{1}{\sqrt{n}} C(\mathcal{E})^n (1 + o(1))$$

as  $n \rightarrow \infty$  for any bounded function  $\phi$  on  $[0, \infty[$  which is continuous a.e. Lebesgue and such that  $\sup_{x>0} |\phi(x)/x| < \infty$ . Or equivalently

$$\int \rho\left(\frac{dP_2^n}{dP_1^n + dP_2^n}\right) d(P_1^n + P_2^n) = \int_0^1 \rho(x) \frac{(1+x)^{t_0}}{x^{1+t_0}} dx \frac{1}{\sqrt{2\pi\tau^2}} \frac{1}{\sqrt{n}} C(\mathcal{E})^n (1 + o(1))$$

as  $n \rightarrow \infty$  for any function  $\rho$  on  $[0, 1]$  which is continuous a.e. Lebesgue and such that  $\sup_{1>x>0} |\rho(x)/x| < \infty$  and  $\sup_{0<x<1} |\rho(x)/(1-x)| < \infty$ .

If  $\psi$  is sublinear or superlinear on  $R^2$  then  $\rho(x) = \psi(1-x, x); x \in [0, 1]$  satisfies the requirements of the Theorem. Thus, by specializing to functions  $x \rightsquigarrow \wedge_t \sum_{\theta=1}^2 L_\theta(t) x_\theta$ , we find asymptotic expressions for minimum Bayes risk in various decision problems.

**4. Replicated experiments when the parameter set is finite.** How fast does the content of information in  $n$  replicates of an experiment  $\mathcal{E}$  increase as  $n \uparrow \infty$ ? In this section we shall investigate this question when  $\Theta$  is finite. In view of the fact that pairwise sufficiency implies sufficiency, it is not too surprising that, up to a first approximation, the problem may be reduced to the same problem for dichotomies. A few crude, but for our purposes sufficient, inequalities are collected in

PROPOSITION 4.1.

(i) Suppose  $\lambda(\Theta_0) > 0$  and let  $\lambda^0$  be the conditional distribution on  $\Theta_0$  given “ $\theta \in \Theta_0$ .” Then

$$\lambda(\Theta_0) b(\lambda^0 | \mathcal{E}_{\Theta_0}) \leq b(\lambda | \mathcal{E}).$$

(ii) If  $\lambda$  is nondegenerate then

$$2b(\lambda | \mathcal{E}) \leq \sum_{\theta_1+\theta_2} (\lambda_{\theta_1} + \lambda_{\theta_2}) b\left(\frac{\lambda_{\theta_1}}{\lambda_{\theta_1} + \lambda_{\theta_2}}, \frac{\lambda_{\theta_2}}{\lambda_{\theta_2} + \lambda_{\theta_1}} \middle| \mathcal{E}_{(\theta_1, \theta_2)}\right).$$

(iii) Suppose  $\Theta$  has  $m < \infty$  elements. Then

$$\delta_a(\mathcal{E}) \leq \min\left\{2mb\left(\frac{1}{m}, \dots, \frac{1}{m} \middle| \mathcal{E}\right), (m-1)\max_{\theta_1+\theta_2} \delta_a(\mathcal{E}_{(\theta_1, \theta_2)})\right\}.$$

PROOF. (i) and part of (iii) follows from the inequalities  $\sum_{\theta_0} \lambda(\theta) r(\theta) \leq \sum \lambda(\theta) r(\theta) \leq \sum r(\theta)$  which are valid for any nonnegative risk function  $r$ . (ii) follows from the inequality

$$\sum \lambda_\theta f_\theta - \vee \lambda_\theta f_\theta \leq \frac{1}{2} \sum_{\theta_1+\theta_2} \lambda_{\theta_1} f_{\theta_2} \wedge \lambda_{\theta_2} f_{\theta_1}$$

valid for any nonnegative numbers  $f_\theta, \theta \in \Theta$ . (iii) follows now from (ii) and (2.9).  $\square$

It will be assumed throughout this section that the parameter set  $\Theta$  is finite. The  $\theta$ th unit vector  $e^\theta$  in  $R^\Theta$  is defined by  $e^\theta(\theta') = 1$  or  $0$  as  $\theta' = \theta$  or  $\theta' \neq \theta$ . We extend the definition of the constant  $C(\mathcal{E})$  in Section 3 by defining

$$C(\mathcal{E}) = \max_{\theta_1+\theta_2} \inf_{0<t<1} \int dP_{\theta_1}^{1-t} dP_{\theta_2}^t.$$

Thus

$$C(\mathcal{E}) = \max_{\theta_1+\theta_2} C(\mathcal{E}_{(\theta_1, \theta_2)}).$$

Consider now an experiment  $\mathcal{E} = (P_\theta; \theta \in \Theta)$ . Let  $\psi$  be sublinear on  $R^\Theta$  and let  $F$  be a nonempty subset of  $\Theta$ . Then, by sublinearity  $\psi(z) = \psi(\sum z_\theta e^\theta) \leq \psi(\sum_F z_\theta e^\theta) + \sum_{F^c} z_\theta \psi(e^\theta)$ . Let  $S$  denote the standard measure of  $\mathcal{E}$ . Then

$$\psi(\mathcal{M}_a) - \psi(\mathcal{E}) = \sum_\theta \psi(e^\theta) - \int \psi dS = \int [\sum z_\theta \psi(e^\theta) - \psi(\sum z_\theta e^\theta)] S(dz)$$

$$\geq \int [\sum_F z_\theta \psi(e^\theta) - \psi(\sum_F z_\theta e^\theta)] S(dz) = \psi_F(\mathcal{M}_\alpha) - \psi_F(\mathcal{E})$$

where  $\psi_F(z) = \psi(\sum_F z_\theta e^\theta)$ . If in particular,  $F = \{\theta_1, \theta_2\}$  then  $\psi(\mathcal{M}_\alpha) - \psi(\mathcal{E}) \geq \psi_{\{\theta_1, \theta_2\}}(\mathcal{M}_\alpha) - \psi_{\{\theta_1, \theta_2\}}(\mathcal{E})$ . Substituting  $\mathcal{E}^n$  for  $\mathcal{E}$  and applying Theorem 3.9 we find, provided  $\psi_{\{\theta_1, \theta_2\}}$  is not affine on  $[0, \infty[^\Theta$ , that  $\liminf_n [\psi(\mathcal{M}_\alpha) - \psi(\mathcal{E}^n)]^{1/n} \geq C(\mathcal{E}_{\{\theta_1, \theta_2\}})$ . Suppose now that this provision is satisfied for all two point sets  $\{\theta_1, \theta_2\}$ . Then

$$\liminf_n [\psi(\mathcal{M}_\alpha) - \psi(\mathcal{E}^n)]^{1/n} \geq C(\mathcal{E}).$$

The provision above is obviously satisfied for any function  $z \rightsquigarrow \vee_\theta \lambda_\theta z_\theta$  where  $\lambda$  is a prior distribution on  $\Theta$  such that  $\lambda_\theta > 0$  for all  $\theta$ . Suppose  $\lambda$  satisfies this condition. Then the above results imply that

$$\liminf_n b(\lambda | \mathcal{E}^n)^{1/n} \geq C(\mathcal{E}).$$

By Proposition 4.1 and Theorem 3.5

$$\limsup_n b(\lambda | \mathcal{E}^n)^{1/n} \leq \max_{\theta_1 + \theta_2} C(\mathcal{E}_{\{\theta_1, \theta_2\}}) = C(\mathcal{E}).$$

It follows that  $b(\lambda | \mathcal{E}^n)^{1/n} \rightarrow C(\mathcal{E})$  as  $n \rightarrow \infty$ . Hence, by (2.9) and Proposition 4.1

$$\begin{aligned} C(\mathcal{E}) &= \lim_n b(\text{uniform} | \mathcal{E}^n)^{1/n} \leq \liminf_n [\delta_\alpha(\mathcal{E}^n)/2]^{1/n} \\ &\leq \limsup_n [\delta_\alpha(\mathcal{E}^n)/2]^{1/n} \leq \limsup_n [(\neq\Theta)b(\text{uniform} | \mathcal{E}^n)]^{1/n} = C(\mathcal{E}), \end{aligned}$$

so that

$$\lim_n \delta_\alpha(\mathcal{E}^n)^{1/n} = C(\mathcal{E}).$$

By the sublinear function criterion

$$\begin{aligned} \limsup_n [\psi(\mathcal{M}_\alpha) - \psi(\mathcal{E}^n)]^{1/n} &\leq \limsup_n [\sum_\theta 1/2 [\psi(e^\theta) + \psi(-e^\theta)] \delta_\alpha(\mathcal{E}^n)]^{1/n} \\ &\leq \limsup_n \delta_\alpha(\mathcal{E}^n)^{1/n} = C(\mathcal{E}) \end{aligned}$$

for any sublinear function  $\psi$  on  $R^\Theta$ . Altogether we have proved the following.

**THEOREM 4.2.** *Let  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  be an experiment with finite parameter set. Then*

- (i)  $\lim_{n \rightarrow \infty} \delta_\alpha(\mathcal{E}^n)^{1/n} = C(\mathcal{E})$ .
- (ii)  $\lim_{n \rightarrow \infty} b(\lambda | \mathcal{E}^n)^{1/n} = C(\mathcal{E})$  provided  $\lambda_\theta > 0$  for all  $\theta$ .
- (iii)  $\limsup_{n \rightarrow \infty} [\psi(\mathcal{M}_\alpha) - \psi(\mathcal{E}^n)]^{1/n} \leq C(\mathcal{E})$  for any sublinear function  $\psi$ .
- (iv)  $\lim_{n \rightarrow \infty} [\psi(\mathcal{M}_\alpha) - \psi(\mathcal{E}^n)]^{1/n} = C(\mathcal{E})$  for any sublinear function  $\psi$  on  $R^\Theta$  such that none of the maps  $z \rightsquigarrow \psi(z_{\theta_1} e^{\theta_1} + z_{\theta_2} e^{\theta_2})$ ;  $\theta_1 \neq \theta_2$  are linear on  $[0, \infty[^\Theta$ .

**REMARK.** Using the inequalities mentioned after (2.11) it may be shown that  $[2 - 2/m - \delta_i(\mathcal{E}^n)]^{1/n} \rightarrow C(\mathcal{E})$ .

If  $\psi(x) = \vee_t \sum_\theta \lambda_\theta U_\theta(t) x_\theta$  where  $T$  is a decision space and  $U$  is the utility function, then (iv) describes the exponential rate of convergence to  $\sum_\theta \lambda_\theta \vee_\theta U_\theta(t)$  of maximum Bayes utility. If the  $\vee_t$  in the expression for  $\psi$  is attained and  $\lambda_\theta > 0$  when  $\theta \in \Theta$ , then the exceptional case is precisely the situation where for some two point set  $\{\theta_1, \theta_2\}$  no observations are needed when it is known that  $\theta \in \{\theta_1, \theta_2\}$ . Thus (iv) is not applicable to expressions like  $\|\wedge_\theta \lambda_\theta P_\theta^n\| = \|\vee_\theta (-\lambda_\theta) P_\theta^n\|$  when  $\neq\Theta \geq 3$  and, in fact,  $\|\wedge_\theta \lambda_\theta P_\theta^n\| \leq \prod \lambda_\theta^n H_\theta(t)^n$  for any pair  $(\lambda, t)$  of prior distributions on  $\Theta$ .

Although Theorem 4.2 yields the exact rate of exponential convergence in many situations, there are situations of interest where the condition in (iv) is not satisfied. Consider, for example, the problem of catching  $\theta$  with an  $r$ -point confidence set. Then the minimax probability of not covering the true value is

$$\kappa_r(\mathcal{E}) = 1 - \inf_\lambda \|\vee_U \sum_{\theta \in U} \lambda_\theta P_\theta\|,$$

where  $U$  runs through all  $r$ -point subsets of  $\Theta$ . Clearly  $\kappa_1 = \delta(\mathcal{E}, \mathcal{M}_a)/2$  and it is easily seen that  $\kappa_r$  is monotonically decreasing in  $r$ .

The minimum Bayes probability of not covering  $\theta$  for the prior  $\lambda$  is  $1 - \|\vee_U \sum_{\theta \in U} \lambda_\theta P_\theta\| = \|\wedge_U \sum_{\theta \in U} \lambda_\theta P_\theta\|$ . This minimum is achieved by the confidence set  $x \rightsquigarrow U_x$  if and only if the class of sets  $U$  such that  $\sum_{\theta \in U} \lambda_\theta f_\theta(x) = \vee_U \sum_{\theta \in U} \lambda_\theta f_\theta(x)$  has the probability 1 for  $\sum_\theta \lambda_\theta P_\theta$  almost all  $x$ . Here  $f_\theta = dP_\theta/d\sum_\theta P_\theta$ ,  $\theta \in \Theta$ .

Let us briefly consider the asymptotic behavior of these quantities.

**THEOREM 4.3.** *Suppose  $\Theta$  is finite and put  $m = \#\Theta$ . Define for each experiment and each integer  $r \in \{1, 2, \dots, m - 1\}$  the quantity  $\kappa_r(\mathcal{E})$  as above. Then*

$$\lim_{n \rightarrow \infty} \kappa_r(\mathcal{E}^n)^{1/n} = \max_W \inf_{t \in \Lambda_W} \int \prod_W (dP_\theta)^{t_\theta}$$

where

- (i)  $W$  runs through all  $(r + 1)$ -point subset of  $\Theta$  and
- (ii)  $\Lambda_W$ , for each  $W$ , is the set of all prior distributions on  $\Theta$  which are supported by  $W$ .

Furthermore, the  $n$ th root of the minimum Bayes probability in  $\mathcal{E}^n$  of not covering  $\theta$  for the prior  $\lambda$  converges to the same limit, provided  $\lambda_\theta > 0$  for all  $\theta \in \Theta$ .

**REMARK.** Suppose the prior distribution  $\lambda$  assigns positive mass to each  $\theta \in \Theta$ . Then any sequence  $U_1, U_2, \dots$  of confidence sets such that  $U_n$ , for each  $n$ , is an optimal Bayes procedure based on  $\mathcal{E}^n$ , achieves the optimal rate.

Putting  $m = 2$  and  $r = 1$  we see that the last statement generalizes Chernoff's result, Theorem 3.7. The proof of Theorem 4.3 is based on Theorem 4.4 and Proposition 4.5 below.

**THEOREM 4.4.** *Let  $\mu_\theta, \theta \in \Theta$  be a finite family of finite nonnegative measures on a common measurable space. Then*

$$\lim_{n \rightarrow \infty} \|\wedge_\theta \mu_\theta^n\|^{1/n} = \sup_n \|\wedge_\theta \mu_\theta^n\|^{1/n} = \inf_{t \in \Lambda} \int \prod_\theta (d\mu_\theta)^{t_\theta}$$

**REMARK.** The statement of the Theorem is clearly equivalent to the following.

Let  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  be an experiment with finite parameter set  $\Theta$  and let  $a_\theta, \theta \in \Theta$  be a family of nonnegative numbers. Then

$$\lim_{n \rightarrow \infty} \|\wedge_\theta a_\theta^n P_\theta^n\|^{1/n} = \sup_n \|\wedge_\theta a_\theta^n P_\theta^n\|^{1/n} = \inf_{t \in \Lambda} \prod_\theta a_\theta^{t_\theta} H(t | \mathcal{E}).$$

**PROOF.** As  $\wedge_\theta f_\theta \leq \prod_\theta f_\theta^{t_\theta}$ ,  $f \in [0, \infty]^\Theta$ ,  $t \in \Lambda$  we get

$$\|\wedge_\theta \mu_\theta^n\| \leq \int \prod_\theta d(\mu_\theta^n)^{t_\theta} = \left( \int \prod_\theta d\mu_\theta^{t_\theta} \right)^n.$$

It remains therefore to show that  $\liminf_n \|\wedge_\theta \mu_\theta^n\|^{1/n} \geq \inf_t \int \prod_\theta d\mu_\theta^{t_\theta}$ . Note that both sides of this inequality remains unchanged if each  $\mu_\theta$  is replaced by its  $\wedge_\theta \mu_\theta$  absolutely continuous component. This reduces the problem to showing that

$$(4.1) \quad \liminf \|\wedge_{i=1}^m e^{-na_i} P_i^n\|^{1/n} \geq \inf_t e^{-(a,t)} H(t | P_1, \dots, P_m)$$

for any homogeneous experiment  $\mathcal{E} = (P_1, \dots, P_m)$  and any  $a \in R^m$  such that  $a_1 = 0$ . If  $m = 2$  then (4.1) is (3.1) in our proof of Chernoff's result in the previous section and the

general case may be proved quite similarly. Again we may assume that the inf is obtained at a  $t^0 \in \Lambda^0$  and we may also assume that  $\mathcal{L}(\log(dP_i/dP_1), i = 2, \dots, m | P_1)$  is nonsingular. By considering the Hellinger transforms it follows again that  $(\mathcal{E}^n: t^0 + s'/\sqrt{n}, \dots, t^0 + s^m/\sqrt{n})$  converges to an experiment  $\mathcal{G} = (Q_1, \dots, Q_m)$ , where (see Example 4.6)  $Q_i, i = 1, 2, \dots, m$ , are the  $(m - 1)$ -variate multinormal distributions with means  $(s_2^i, s_3^i, \dots, s_m^i)$  and covariance matrix  $\sigma^{-1}$  given by  $\sigma_{\alpha\beta} = \rho_0^{-1} \frac{\partial^2}{\partial t_\alpha \partial t_\beta} \rho |_{t=t_0}, \rho(t) = e^{-(a,t)} H(t | \mathcal{E})$  and  $\rho_0 = \inf_t \rho(t)$ . If  $U_n = dP_n^n/dP_1^n$  and  $U_i = dQ_i/dQ_1$  then this implies that

$$E_{P_1^n} \phi(U_n^{1/\sqrt{n}} e^{-\sqrt{n}a_i}; i = 2, \dots, m) \prod_{i=2}^m U_{ni}^0 \sim \rho_0^n e^{na_i^0} E_{Q_1} \phi(e^{1/2\sigma_{ii}} U_i, i = 2, \dots, m)$$

when  $\phi$  is bounded and continuous a.e. If  $\phi(x_2, \dots, x_m) = 1$  or  $0$  as  $(x_2, \dots, x_m) \in [M^{-1}, M]^m$  or not then we get  $\| \wedge_i e^{-na_i} P_i^n \| \geq \gamma_n$  where  $\gamma_n \sim M^{-\sqrt{n}(2-t_0^i)} \rho_0^n E_{Q_1} \phi(e^{\sigma_{ii}/2} U_i, i = 2, \dots, m)$  and this yield the desired inequality.  $\square$

PROPOSITION 4.5. For any experiment  $\mathcal{E} = (P_1, P_2, \dots, P_m)$ , if  $1 \leq r < m$  then

$$\sum_{B_{r+1}} \binom{m-1}{r}^{-1} \kappa_r(\mathcal{E}_{\{i_1, i_2, \dots, i_{r+1}\}}) \leq \kappa_r(\mathcal{E}) \leq \sum_{B_{r+1}} \kappa_r(\mathcal{E}_{\{i_1, i_2, \dots, i_{r+1}\}}),$$

where  $B_{r+1} = \{(i_1, \dots, i_{r+1}): i_1 < i_2 < \dots < i_{r+1}\}$ .

PROOF. We may write  $\kappa_r(\mathcal{E}) = \sup_{\lambda} \| \wedge_U \sum_{U'} \lambda_\theta P_\theta \|$ . Let  $y_1 \leq y_2 \leq y_3 \leq \dots \leq y_m$ . Then  $\wedge_U \sum_{U'} y_\theta = y_1 + \dots + y_s$ , where  $s = m - r$ . The inequalities follow now immediately from the identity  $\sum \{y_{i_1} : i_1 < i_2 < \dots < i_{r+1}\} = \binom{m-1}{r} y_1 + \binom{m-2}{r} y_2 + \dots + \binom{r}{r} y_s$ .  $\square$

PROOF OF THEOREM 4.3. The Theorem follows by applying Proposition 4.5 to  $\mathcal{E}^n$ , taking the  $n$ th root, letting  $n \rightarrow \infty$  and using Theorem 4.4.  $\square$

EXAMPLE 4.6. Let  $P_i, i = 1, \dots, m$ , be the  $p$ -variate normal distributions with means  $\xi^i$  and nonsingular covariance matrix  $M$ . Then  $H(t | P_1, \dots, P_m) = e^{-1/4Q}$  where  $Q = \sum t_i t_j (\xi^i - \xi^j)' M^{-1} (\xi^i - \xi^j)$ . If, in particular,  $M$  is the identity matrix and  $\xi^1, \dots, \xi^m$  is an orthonormal basis for  $R^m$  then

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \kappa_r(P_1^n, P_2^n, \dots, P_m^n) = -\frac{r}{2(r+1)}.$$

If the map  $\theta \rightsquigarrow P_\theta$  is not 1 - 1 then  $C(\mathcal{E}) = 1$  and  $\delta(\mathcal{E}^n, \mathcal{M}_\alpha) \geq 1$  for all  $n$ . The obvious way out is to replace the parameter set  $\Theta$  by the set  $\{P_\theta; \theta \in \Theta\}$ . We omit the details.

**5. The general case.** It follows, since deficiencies decrease by taking restrictions, that  $\liminf_n \delta_\alpha(\mathcal{E}^n)^{1/n} \geq C(\mathcal{E})$  for any experiment  $\mathcal{E} = (P_\theta; \theta \in \Theta)$  where  $C(\mathcal{E}) = \sup_{\theta_1, \theta_2} C(\mathcal{E}_{\{\theta_1, \theta_2\}})$ . Although providing a lower bound,  $C(\mathcal{E})$  alone does not determine the exponential rate of convergence when  $\Theta$  is infinite. We shall, nevertheless, see that  $\delta_\alpha(\mathcal{E}^n)^{1/n}$  always converges as  $n \rightarrow \infty$  to the quantity  $\sigma(\mathcal{E}) = \max(C(\mathcal{E}), \tau(\mathcal{E}))$  where  $\tau(\mathcal{E}) = \inf_n (\delta_\alpha(\mathcal{E}^n)/2)^{1/n}$ . In order to establish this and related results we shall have to search for estimators  $\hat{\theta}_n$  making  $P_\theta^n(\hat{\theta}_n \neq \theta)$  small.

Suppose that an estimator  $\hat{\theta} = \hat{\theta}(X)$  of  $\theta$  based on one observation  $X$  of  $\mathcal{E} = (\chi, \mathcal{A}, P_\theta; \theta \in \Theta)$  is available. Let  $X_1, \dots, X_n$  be  $n$  independent observations of  $X$  and put  $\theta_i = \hat{\theta}(X_i); i = 1, 2, \dots, n$ . Restrict the search to the subset  $U$  of  $\Theta$  consisting of those  $\theta'$  whose relative frequency  $h_n(\theta')$  in  $(\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_n)$  exceed a fixed number  $\xi \in ]0, 1]$ . Then pick an estimator  $s_n = s_n(X_1, \dots, X_n)$  within the subset of  $U$  consisting of those  $\theta'$  which maximizes the likelihood *within*  $U$ . Thus  $s_n$  is a restricted maximum likelihood estimator.

To be more precise we may proceed as follows. Equip  $\Theta$  with the  $\sigma$ -algebra generated

by the one-point sets and suppose  $\theta$  is measurable. Put  $g_{\theta, F} = dP_{\theta}^n/d\sum_{\theta \in F} P_{\theta}^n$  whenever  $F$  is a finite subset of  $\Theta$  containing  $\theta$ . Choose for each nonempty finite subset  $F$  of  $\Theta$  a measurable estimator  $t_F = t_F(X_1, \dots, X_n) \in F$  such that  $g_{t_F, F} = \vee_{\theta \in F} g_{\theta, F}$ . We may then take  $s_n = t_U$  provided  $U \neq \emptyset$ . Put  $s_n = \theta_0$  where  $\theta_0$  is a fixed element of  $\Theta$  when  $U = \emptyset$ . Then  $s_n$  is measurable with respect to the completion of  $\mathcal{A}^n$  for  $P_{\theta}^n$  for each  $\theta$  and

**PROPOSITION 5.1.** Put  $r = r(\xi) = \xi^{-\xi}(1 - \xi)^{\xi-1}$ ,  $C_{\theta} = \sup\{C(\mathcal{E}_{\{\theta, \theta'\}}; \theta' \neq \theta)\}$  and assume  $0 < \eta \leq \xi \leq P_{\theta}(\hat{\theta} = \theta)$ . Then

$$P_{\theta}^n(s_n \neq \theta) \leq [P_{\theta}(\hat{\theta} \neq \theta)^{1-\xi} r]^n + \eta^{-1} [r\eta^{\xi}]^n + \eta^{-1} C_{\theta}^n.$$

**PROOF.** Put  $F_{\theta} = \{\theta' : P_{\theta}(\hat{\theta} = \theta') \geq \eta\}$ . Then

$$P_{\theta}^n(s_n \neq \theta) \leq P_{\theta}^n(U \not\subseteq F_{\theta}) + P_{\theta}^n(\theta \notin U) + P_{\theta}^n(\theta \in U \not\subseteq F_{\theta} \text{ \& } s_n \neq \theta).$$

Now

$$P_{\theta}^n(U \not\subseteq F_{\theta}) = \sum' P_{\theta}^n(U \not\supseteq \theta') = \sum' P_{\theta}^n(h_n(\theta') \geq \xi)$$

where  $\sum' = \sum_{\theta' \notin F_{\theta}}$ . Applying the inequality  $P_{\theta}(Z \geq C) \leq \inf_{t \geq 0} e^{-tC} E_{\theta} e^{tZ}$  to each term we find

$$P_{\theta}^n(U \not\subseteq F_{\theta}) \leq \sum' [P_{\theta}(\hat{\theta} \neq \theta')^{1-\xi} P_{\theta}(\hat{\theta} = \theta')^{\xi} r]^n \leq \sum' P_{\theta}(\hat{\theta} = \theta') P_{\theta}(\hat{\theta} = \theta')^{n\xi-1} r^n \leq \eta^{n\xi-1} r^n$$

provided  $n\xi \geq 1$ . Hence  $P_{\theta}^n(U \not\subseteq F_{\theta}) \leq \eta^{n\xi-1} r^n$  in any case. Similarly

$$P_{\theta}^n(\theta \notin U) = P_{\theta}^n(h_n(\theta) < \xi) \leq [P_{\theta}(\hat{\theta} \neq \theta)^{1-\xi} P_{\theta}(\hat{\theta} = \theta)^{\xi} r]^n \leq [P_{\theta}(\hat{\theta} \neq \theta)^{1-\xi} r]^n.$$

Finally, by the remark to Theorem 3.8, we find

$$P_{\theta}^n(\theta \in U \not\subseteq F_{\theta}, s_n \neq \theta) = \sum_{\theta' \in F_{\theta} - \{\theta\}} P_{\theta}^n(\theta \in U \not\subseteq F_{\theta} \text{ \& } s_n = \theta')$$

$$\leq \sum_{\theta' \in F_{\theta} - \{\theta\}} P_{\theta}^n(g_{\theta', F_{\theta}} \geq g_{\theta, F_{\theta}}) \leq \sum_{\theta' \in F_{\theta} - \{\theta\}} \|P_{\theta'}^n \wedge P_{\theta}^n\| \leq \#(F_{\theta} - \{\theta\}) C_{\theta}^n \leq \eta^{-1} C_{\theta}^n. \quad \square$$

**COROLLARY 5.2.** Suppose  $\mathcal{E}^k$  is  $\varepsilon = (\varepsilon_{\theta}; \theta \in \Theta)$  deficient with respect to  $\mathcal{M}_{\alpha}$  and let  $0 < \eta \leq \xi \leq 1$ . Then  $\mathcal{E}^n$  is  $(\gamma_{n\theta}; \theta \in \Theta)$  deficient with respect to  $\mathcal{M}_{\alpha}$  where

$$\frac{1}{2} \gamma_{n\theta} = [(\varepsilon_{\theta}/2)^{1-\xi} r]^{[n/k]} + \eta^{-1} [r\eta^{\xi}]^{[n/k]} + \eta^{-1} C_{\theta}^{[n/k]}$$

or  $1$  as  $\xi_{\theta} \leq 1 - (\varepsilon_{\theta}/2)$  or not.

**REMARK.** Choosing  $\eta = 1/n$  and  $\xi = (\log n)^{-1/2}$  we find that there is a sequence  $s_1, s_2, \dots$  of estimators such that

$$\limsup_n P_{\theta}^n(\hat{s}_n \neq \theta)^{1/n} \leq (\varepsilon_{\theta}/2)^{1/k} \vee C_{\theta}.$$

Here is the main result of this section.

**THEOREM 5.3.**

$$\delta_{\alpha}(\mathcal{E}^n)^{1/n} \rightarrow \sigma(\mathcal{E}) \quad \text{as } n \rightarrow \infty.$$

**REMARK.** It follows that  $\liminf_n \sup_{\theta} P_{\theta}^n(s_n \neq \theta)^{1/n} \geq \sigma(\mathcal{E})$  for any sequence  $s_n(X_1, \dots, X_n); n = 1, 2, \dots$  of estimators of  $\theta$  and that  $\lim_n \sup_{\theta} P_{\theta}^n(s_n \neq \theta)^{1/n} = \sigma(\mathcal{E})$  for some sequence  $(s_1, s_2, \dots)$  of estimators.

**PROOF.** Clearly  $\liminf_n \delta_{\alpha}(\mathcal{E}^n)^{1/n} \geq \sigma(\mathcal{E})$ . It remains to show that  $\limsup_n \delta_{\alpha}(\mathcal{E}^n)^{1/n}$



$\leq C(\mathcal{E}) \vee (\varepsilon_k/2)^{1/k}$  when  $\varepsilon = \delta_\alpha(\mathcal{E}^k) < 2$ . This, however, follows from the Corollary by letting  $\eta \rightarrow 0$  and then  $\xi \rightarrow 0$  in the inequality  $\limsup_n \gamma_n^{1/n} \leq ((\varepsilon/2)^{1-\xi r})^{1/k} \vee (r\eta^\xi)^{1/k} \vee C(\mathcal{E})$ ,  $\eta < \xi < 1 - \varepsilon/2$ .  $\square$

Applying the Theorem to  $\mathcal{E}^k$  we find easily that  $\sigma(\mathcal{E}^k) = \sigma(\mathcal{E})^k$  and clearly  $\sigma(\mathcal{E}) \cong \sigma(\mathcal{F})$  when  $\mathcal{E} \cong \mathcal{F}$ .

**COROLLARY 5.4.** *The following conditions are equivalent for an experiment  $\mathcal{E} = (P_\theta; \theta \in \Theta)$ ,*

- (i)  $\delta_\alpha(\mathcal{E}^n) < 1$  for some  $n$ ;
- (ii)  $\lim_n \delta_\alpha(\mathcal{E}^n) = 0$ ;
- (iii)  $\delta_\alpha(\mathcal{E}^n) \leq c\rho^n$ ,  $n = 1, 2, \dots$  for some constant  $c > 0$  and some constant  $\rho < 1$ ;
- (iv)  $\delta_\alpha(\mathcal{E}^n) < 2$  for some  $n$  and  $\inf_{\theta_1, \theta_2} \|P_{\theta_1} - P_{\theta_2}\| > 0$ .

**PROOF.** Suppose  $\delta_\alpha(\mathcal{E}^r) < 1$ . Then  $\tau(\mathcal{E}) \leq \delta_\alpha(\mathcal{E}^r)^{1/r} < 1$  and  $\sup_{\theta_1, \theta_2} \delta_\alpha(\mathcal{E}^r_{\{\theta_1, \theta_2\}}) < 1$ . The last inequality is equivalent to  $\inf_{\theta_1, \theta_2} \delta_i(\mathcal{E}^r_{\{\theta_1, \theta_2\}}) > 0$  or  $C(\mathcal{E})^r = \sup_{\theta_1, \theta_2} C(\mathcal{E}^r_{\{\theta_1, \theta_2\}}) < 1$ . Hence, as  $\tau(\mathcal{E}) < 1$  and  $C(\mathcal{E}) < 1$ ,  $\sigma(\mathcal{E}) < 1$ . Thus, by the Theorem, (i)  $\Rightarrow$  (iii). The other implications of the Corollary are then straightforward.  $\square$

It follows that if  $\mathcal{E}^n \rightarrow \mathcal{M}_\alpha$  then the speed of convergence is necessarily exponential provided  $\mathcal{E} \sim \mathcal{M}_\alpha$ .

The constant 1 in (i) cannot be increased. If, for example,  $\mathcal{E} = (P_1, P_2, P_3, \dots)$ , where  $P_1 = P_2$  and  $P_i \wedge P_j = 0$  when  $i, j \geq 2$ , and  $\delta_\alpha(\mathcal{E}^n) = 1$  for all  $n = 1, 2, \dots$ .

We saw in the previous section that  $C(\mathcal{E})$  alone determined the rate of convergence when  $\Theta$  is finite. Here is an example, showing that this (i.e.,  $C(\mathcal{E}) \cong \tau(\mathcal{E})$ ) does not hold in general when  $\Theta$  is infinite.

**EXAMPLE 5.5.** Suppose  $\Theta = \{1, 2, \dots\}$  and that the density of  $P_\theta$  with respect to the uniform distribution  $P$  on  $[0, 1]$  is  $f_\theta$  where  $f_\theta(x) = 2$  or  $= 0$  according to whether  $x$  belongs to one of the intervals  $[(k-1)/2^\theta, (k-1/2)/2^\theta]; k = 1, \dots, 2^\theta$  or not. Then  $E_P f_\theta(X)^{1-t} f_{\theta_2}(X)^t = 1/2$ ,  $\theta_1 \neq \theta_2$ ,  $0 < t < 1$ . It follows that the dichotomies  $\mathcal{E}_{\{\theta_1, \theta_2\}}; \theta_1 \neq \theta_2$  are all equivalent to the simple dichotomy  $((1/2, 0, 1/2), (0, 1/2, 1/2))$  and, in particular, that  $C(\mathcal{E}) = 1/2$ . On the other hand it is not difficult to see that  $P_\theta(B) \rightarrow P(B)$  as  $\theta \rightarrow \infty$  for any Borel set  $B$ . By the Proposition below  $\delta_\alpha(\mathcal{E}^n) \equiv_n 2$ .

**PROPOSITION 5.6.**  $\delta_\alpha(\xi^n) \equiv_n 2$  provided either

- (i)  $\Theta$  contains an uncountable subset  $\Theta_0$  such that  $(P_\theta; \theta \in \Theta_0)$  is dominated, or
- (ii)  $\Theta$  contains distinct elements  $\theta_1, \theta_2, \dots$  such that  $\lim_n P_{\theta_n}(A)$  exists for all  $A \in \mathcal{A}$ .

**REMARK.** Let  $\mu$  be a probability measure dominating  $\mathcal{E}$  and put  $f_\theta = dP_\theta/d\mu$ . Then condition (ii) is satisfied whenever  $\Theta$  is infinite and  $(f_\theta, \theta \in \Theta)$  are uniformly integrable. Thus  $\delta_\alpha(\mathcal{E}^n) \equiv_n 2$  when the sample space of  $\mathcal{E}$  is finite and  $\Theta$  is infinite. In spite of this  $\mathcal{E}^n$  converges weakly, i.e., for restrictions to finite subparameter sets, whenever  $\theta \rightsquigarrow P_\theta$  is 1 - 1.

**PROOF.** The sufficiency of condition (i) was noted after (2.7). We may then without loss of generality assume that  $\Theta = \{1, 2, \dots\}$  and that  $P_n(A) \rightarrow P(A)$  for any event  $A$ . If  $M$  is an estimator of  $\theta$  then  $P_n(M = n) \leq P_n(M \geq N) \rightarrow 0$  as  $n \rightarrow \infty$  and then  $N \rightarrow \infty$ . Hence  $\inf_\theta P_\theta(M = \theta) = 0$  so that  $\delta_\alpha(\mathcal{E}) = 2$ . The Proposition follows, since by the Vitali-Hahn-Saks Theorem,  $P'_n \rightarrow P'$  in the same sense.  $\square$

If  $\Theta$  is finite and  $\delta_\alpha(\mathcal{E})$  equals its maximal value then  $\mathcal{E} \sim \mathcal{M}_i$  and replications do not yield any information. In the infinite case, however, such a discouraging start does not (see

Example 5.9) prevent information from being accumulated rapidly and with arbitrarily small  $\sigma(\mathcal{E})$ , as the number of replications increase.

How do distances and deficiencies for replicated experiments behave in general? It is clear by Theorem 5.3 and the inequality  $\delta_a(\mathcal{E}^n) - \delta_a(\mathcal{F}^n) \leq \delta(\mathcal{E}^n, \mathcal{F}^n) \leq \delta_a(\mathcal{E}^n)$  that  $\limsup_n \delta(\mathcal{E}^n, \mathcal{F}^n)^{1/n} \leq \sigma(\mathcal{E})$  and that  $\delta(\mathcal{E}^n, \mathcal{F}^n)^{1/n} \rightarrow \sigma(\mathcal{E})$  when  $\sigma(\mathcal{E}) > \sigma(\mathcal{F})$ . The problem of the asymptotic behavior of  $\delta(\mathcal{E}^n, \mathcal{F}^n)$  when  $\sigma(\mathcal{E}) < \sigma(\mathcal{F})$  is open. Our knowledge about the asymptotic behaviour of  $\Delta$  distances is more complete. It follows from the considerations above that  $\limsup_n \Delta(\mathcal{E}^n, \mathcal{F}^n)^{1/n} \leq \sigma(\mathcal{E}) \vee \sigma(\mathcal{F})$  and that  $\Delta(\mathcal{E}^n, \mathcal{F}^n)^{1/n} \rightarrow \sigma(\mathcal{E}) \vee \sigma(\mathcal{F})$  when  $\sigma(\mathcal{E}) \neq \sigma(\mathcal{F})$ . If  $\Theta$  is a two point set then by [26] the above holds whenever  $\mathcal{E} \sim \mathcal{F}$ . Thus it also holds whenever  $\Theta$  is finite and  $\mathcal{E}_{\{\theta_1, \theta_2\}} \sim \mathcal{F}_{\{\theta_1, \theta_2\}}$  for  $\theta_1 \neq \theta_2$ .

EXAMPLE 5.7. Blackwell [2] considered the experiments  $\mathcal{E}$  and  $\mathcal{F}$  given by the matrices

$$\mathcal{E}: \begin{array}{c|cc} \theta \backslash x & 0 & 1 \\ \hline 1 & 1 & 0 \\ 2 & \frac{1}{2} & \frac{1}{2} \\ 3 & 0 & 1 \end{array} \quad \text{and} \quad \mathcal{F}: \begin{array}{c|cc} \theta \backslash x & 0 & 1 \\ \hline 1 & 1 & 0 \\ 2 & \frac{1}{2} & \frac{1}{2} \\ 3 & \frac{1}{2} & \frac{1}{2} \end{array}$$

It is easily seen that  $\mathcal{E}$  is pairwise more informative than  $\mathcal{F}$  and that  $H(t | \mathcal{E}) \leq H(t | \mathcal{F})$  for all  $t$ . Consider, however, decision problems with safety bounds for the risk at  $\theta = 2$ . Then a very good performance in  $\mathcal{E}$  when  $\theta = 3$  will easily lead to a bad performance for  $\theta = 1$  and vice versa. The experiment  $\mathcal{F}$ , on the other hand, is obviously not as bad in this respect.

By sufficiency  $\mathcal{E}^n$  and  $\mathcal{F}^n$  may be reduced respectively to

$$\begin{array}{lcl} x & : & 0, *, n \\ P_{n,1} & : & 1, 0, 0 \\ P_{n,2} & : & \frac{1}{2^n}, 1 - \frac{1}{2^{n-1}}, \frac{1}{2^n} \\ P_{n,3} & : & 0, 0, 1 \end{array} \quad \text{and} \quad \begin{array}{lcl} x & : & 0, n \\ Q_{n,1} & : & 1, 0 \\ Q_{n,2} & : & \frac{1}{2^n}, 1 - \frac{1}{2^n} \\ Q_{n,3} & : & \frac{1}{2^n}, 1 - \frac{1}{2^n} \end{array}$$

Let  $M$  be the randomization from  $\{0, *, n\}$  to  $\{0, n\}$  such that  $M(n | 0) = (2^{2n} + 2^{n+1})^{-1}$ ,  $M(n | *) = 1$  and  $M(n | n) = (2^n + 1 - 2^{-n})(2^n + 2)^{-1}$ . Then  $\|P_{n,i}M - Q_{n,i}\| = 2(2^{2n} + 2^{n+1})^{-1}$ ;  $i = 1, 2, 3$  so that  $\delta(\mathcal{E}^n, \mathcal{F}^n) \leq 2(2^{2n} + 2^{n+1})^{-1}$ .

Now  $(1, 1/2^n, 1/2^n)$  is an available power function in  $\mathcal{F}^n$ . Hence, by testing criterion [22], there is a power function  $\pi$  in  $\mathcal{E}^n$  so that

$$\pi(1) \geq 1 - \frac{\delta}{2}, \quad \pi(2) \leq \frac{1}{2^n} + \frac{\delta}{2} \quad \text{and} \quad \pi(3) \geq \frac{1}{2^n} - \frac{\delta}{2}$$

where  $\delta = \delta(\mathcal{E}^n, \mathcal{F}^n)$ . We may therefore, by sufficiency, write

$$\begin{aligned} \frac{1}{2^n} + \frac{\delta}{2} \geq \pi(2) &= \frac{1}{2^n} (a + c) + \left(1 - \frac{1}{2^{n-1}}\right) b \geq \frac{1}{2^n} (a + c) \\ &= \frac{1}{2^n} (\pi(1) + \pi(3)) \geq \frac{1}{2^n} \left[1 - \frac{\delta}{2} + \frac{1}{2^n} - \frac{\delta}{2}\right]. \end{aligned}$$

where  $a, b, c \in [0, 1]$ . It follows that  $\delta \geq 2(2^{2n} + 2^{n+1})^{-1}$  so that  $\delta(\mathcal{E}^n, \mathcal{F}^n) = 2(2^{2n} + 2^{n+1})^{-1}$ . Hence  $\delta(\mathcal{E}^n, \mathcal{F}^n)^{1/n} \rightarrow 1/4$  while  $C(\mathcal{E}) = 1/2$  and  $C(\mathcal{F}) = 1$ .

EXAMPLE 5.8. (Translation experiments on the integers.) For each distribution  $P$  on the integers and each integer  $\theta$ , let  $P_\theta$  be the right  $\theta$ -translate of  $P$ . Consider the experiment  $\mathcal{E}_P = (P_\theta; \theta \in \Theta)$ . It was shown in [25] that

$$\delta_a(\mathcal{E}^n) = 2(1 - \sum_{x_2, \dots, x_n} \max_n P(x)P(x + x_2) \cdots P(x + x_n)).$$

It follows in particular that  $\delta_a(\mathcal{E}^n) = 2(1 - P(a))^n$  when  $P(a) \geq P(a + 1) \geq \dots$  and  $\sum_{x=a}^\infty P(x) = 1$ .

Furthermore it is not difficult to see that there are translation invariant maximum likelihood estimators  $\hat{\theta}_n$  based on  $\mathcal{E}^n$  and that any such estimator is minimax for the  $0 - 1$  loss estimation problem i.e.,  $2P_\theta(\hat{\theta}_n \neq \theta) \equiv_\theta \delta_a(\mathcal{E}^n)$ . Thus  $\delta_a(\mathcal{E}^n) < 2$  for all  $n$  and clearly  $\inf_{\theta \neq \theta_0} \|P_{\theta_0} - P_\theta\| > 0$ . Hence  $P_\theta(\hat{\theta}_n \neq \theta) = \delta_a(\mathcal{E}^n)/2 \leq C\rho^n$  for some constant  $C > 0$  and for some  $\rho < 1$ .

In spite of this there are translation experiments  $\mathcal{E}_P$  with  $P$  nondegenerate such that two replications are not better than one in the sense of minimax risk. For example, if  $P(x) = p \left( \frac{1-p}{3-p} \right)^{|x|}$ ;  $x = \dots, -1, 0, 1, \dots$  when  $p \in [0, 1[$  then  $X_1$  alone is a translation invariant maximum likelihood estimator for  $\theta$  based on two observations  $(X_1, X_2)$ .

One would expect that if we are able to do extremely well with two observations, then we should be able to do at least moderately well with one observation. If  $\Theta$  is finite then this follows immediately from the compactness of  $\Delta$ -convergence. If the parameter set is infinite, however, then this is not necessarily true. We shall here satisfy ourselves with an example of a translation experiment where we may, on the basis of two observations, guess  $\theta$  with marvelous accuracy while any estimator based on one observation is quite inaccurate. To make things more concrete we may choose the constant so that the probability of a wrong guess for a translation invariant maximum likelihood estimator based on two observations is less than  $10^{-200}$  while, on the other hand, the probability of making a wrong guess is greater than  $1 - 10^{-200}$  for some  $\theta$  for any estimator based on one observation.

EXAMPLE 5.9. Let  $P$  be the uniform distribution on  $\{1, 2, 4, \dots, 2^{N-1}\}$ . Then, by the example in Section 2 in [25],

$$\delta_a(\mathcal{E}_P^n)/2 = P_\theta(\hat{\theta}_n \neq \theta) = (N - 1)N^{-n}, \quad n = 1, 2, \dots$$

where  $\hat{\Theta}_n(X_1, \dots, X_n) = \min \cap_{i=1}^n \{X_i - 1, X_i - 2, X_i - 4, \dots, X_i - 2^{N-1}\}$ . It follows that  $\lim_{N \rightarrow \infty} \delta_a(\mathcal{E}_P^n) = 1$  or  $0$  as  $n = 1$  or  $n \geq 2$ .

If only one observation  $X_1$  is available, then  $\theta$  is located in the  $N$ -point set  $\{1 - X_1, 2 - X_1, \dots, 2^{n-1} - X_1\}$ . If, however, another observation  $X_2$  is available and  $X_2 \neq X_1$ , (this has probability  $1 - 1/N$ ), then  $\theta$  is completely known. Thus we see that the phenomenon is related to the uniqueness of dyadic expansions.

This example may be sharpened by exhibiting a family  $\mathcal{E}_\epsilon$ ;  $\epsilon > 0$  of experiments, necessarily not translation experiments on the integers, such that the deficiency  $\delta_a(\mathcal{E}_\epsilon) \equiv_\epsilon 2$  while  $\delta_a(\mathcal{E}_\epsilon^n) \leq \epsilon^{n-1}$  for all  $\epsilon > 0$  when  $n \geq 2$ .

### REFERENCES

[1] BAHADUR, R. R. (1960). Asymptotic efficiency of tests and estimates. *Sankhyā* **22** 229-252.  
 [2] BLACKWELL, D. (1951). Comparison of experiments. *Proc. Second Berkeley Symposium Math. Statist. Prob.* 93-102.  
 [3] BLACKWELL, D. (1953). Equivalent comparisons of experiments. *Ann. Math. Statist.* **24** 265-272.  
 [4] BLACKWELL, D. and GIRSHICK, M. A. (1954). *Theory of Game and Statistical Decisions*. Wiley, New York.  
 [5] BOHNENBLUST, F., SHAPLEY, L., SHERMANN, S. (1951). Unpublished Rand Memorandum.  
 [6] CHERNOFF, H. (1952). A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Ann. Math. Statist.* **23** 493-507.  
 [7] CHERNOFF, H. (1956). Large sample theory-parametric case. *Ann. Math. Statist.* **27** 1-22.  
 [8] DE GROOT, M. H. (1962). Uncertainty, information and sequential experiments. *Ann. Math. Statist.* **33** 404-419.  
 [9] DVORETZKY, A. WALD and WOLFOWITZ, J. (1951). Elimination of randomization in certain statistical decision procedures and zero-sum two person games. *Ann. Math. Statist.* **22** 1-21.  
 [10] EFRON, B. (1967). The power of the likelihood ratio test. *Ann. Math. Statist.* **38** 802-806.

- [11] EFRON, B. and TRUAX, D. R. (1968). Deviation theory in exponential families. *Ann. Math. Statist.* **39** 1402-1424.
- [12] HELGELAND, J. (1979). Additional observations and statistical information in the case of 1-parameter exponential distribution. *Statist. Res. Report.*, Univ. of Oslo.
- [13] HEYER, H. (1973). *Mathematische Theorie Statistischer Experimente*. Springer, Berlin.
- [14] LE CAM, L. (1964). Sufficiency and approximate sufficiency. *Ann. Math. Statist.* **35** 1419-1455.
- [15] LE CAM, L. (1968). *Théorie Asymptotique de la Décision Statistique*. Les Presses de l'Université de Montréal.
- [16] LE CAM, L. (1974). On the information contained in additional observations. *Ann. Statist.* **2** 630-649.
- [17] LE CAM, L. (1974). Notes on asymptotic methods in statistical decision theory. Centre de Recherches Math., Univ. de Montréal.
- [18] LE CAM, L. (1975). Distances between experiments. In *A Survey of Statistical Design and Linear Models* (J. N. Srivastava, ed.) 383-395. North Holland, Amsterdam.
- [19] LINDQVIST, B. (1977). How fast does a Markov chain forget the initial state? A decision theoretic approach. *Scand. J. Statist.* **4** 145-152.
- [20] LYAPUNOV, A. (1940). Sur les fonctions-vecteurs complètement additives. *Izvestia Akad. Nauk. SSR. Ser. Mat.* **4** 465-478.
- [21] MORSE, N. and SACKSTEDER, R. (1966). Statistical isomorphism. *Ann. Math. Stat.* **37** 203-214.
- [22] SWENSEN, A. R. (1980). Deficiencies in linear normal experiments. *Ann. Statist.* **8** 1142-1155.
- [23] TORGERSEN, E. N. (1970). Comparison of experiments when the parameter space is finite. *Z. Wahrscheinlichkeitstheorie und verw. Gebiete.* **16** 219-249.
- [24] TORGERSEN, E. N. (1972a). Local comparison of experiments. *Statist. Res. Report.*, Univ. of Oslo.
- [25] TORGERSEN, E. N. (1972b). Comparison of translation experiments. *Ann. Math. Statist.* **43** 1383-1399.
- [26] TORGERSEN, E. N. (1974). Asymptotic behaviour of powers of dichotomies. *Statist. Res. Report.*, Univ. of Oslo.
- [27] TORGERSEN, E. N. (1975). *Notes on comparison of statistical experiments*. Chapters 0-8. Recorded and supplemented by B. Lindqvist. *Statist. Memoir.*, Univ. of Oslo.
- [28] TORGERSEN, E. N. (1976a). Comparison of statistical experiments. *Scand. J. Statist.* **3** 186-208.
- [29] TORGERSEN, E. N. (1976b). Deviations from total information and total ignorance as measures of information. *Proc. Fourth Int. Conf. Math. Statist.* Wisla, Poland.
- [30] TORGERSEN, E. N. (1977). Mixtures and products of dominated experiments. *Ann. Statist.* **5** 44-64.

INSTITUTE OF MATHEMATICS  
UNIVERSITY OF OSLO  
P.O. BOX 1053  
BLINDERN, OSLO 3  
NORWAY