

Lectures on Optimal Design and
Sequential Analyses

Herman Chernoff

Inst. of Statistics Mimeo Series # 1349

Preface

Professor Herman Chernoff, Massachusetts Institute of Technology, delivered the Harold Hotelling Lectures 1981, on January 26-29, under the titles

January 26, 1981: Optimal Design of Experiments

January 27, 1981: Sequential Design of Experiments

January 28, 1981: Continuous Time Sequential Problems

January 29, 1981: Massachusetts Number Game

These are the Lecture Notes of the first three lectures, with contributions by Paul P. Gallo and Leonard A. Stefanski. Also taking notes were Ying So, Kenneth Risko, David Giltinan, Jed Frees, Bruce J. Collings and Reuel L. Smith. Paul Gallo and myself did the proofreading of the typescripts.

I. M. Chakravarti

June 1981

LECTURES ON OPTIMAL DESIGN AND SEQUENTIAL ANALYSES

Herman Chernoff
Massachusetts Institute of Technology

I. Introduction

These notes are the framework of the Hotelling Lectures, a series of three lectures presented at the Department of Statistics at the University of North Carolina at Chapel Hill on January 26-28, 1981.

These lectures touch briefly on three topics which are discussed more fully in the SIAM Monograph No. 8, entitled *Sequential Analysis and Optimal Design*. The titles of the lectures are

- 1) Optimal Design of Experiments for Fixed Sample Size,
- 2) Sequential Design of Experiments, and
- 3) Continuous Time Sequential Problems

Many of the ideas and results are introduced in terms of examples with relatively little explanatory introduction, but I believe that these notes present a relatively comprehensible glimpse of a subject of importance and interest.

II. Optimal Design of Experiments

A. Optimal Design in Estimation

We first consider some fixed sample size estimation problems:

Example 1 X_1, X_2, \dots, X_n are i.i.d. $N(\mu, \sigma^2)$, with σ^2 known. We wish to estimate μ ; how large a sample should be taken? Costs involved:

- (i) cost of sampling: it is often reasonable to assume that this is linear in n ;

(ii) "cost" of bad decision: squared error loss is often used.

Thus, if we use the sample mean \bar{X} to estimate μ , then for some c, k , the cost is

$$cn + (\bar{X} - \mu)^2 k$$

We want to minimize the expected cost:

$$cn + k\sigma^2/n$$

This is done by choosing $n = n_0 = (k\sigma^2/c)^{1/2}$ and the optimal cost is $c_0 = 2(ck)^{1/2}\sigma$.

We now consider regression experiments in which we are interested in estimating some specific functions of regression parameters.

Example 2 $Y_i = \alpha + \beta x_i + u_i, i = 1, \dots, n; u_i$ i.i.d. $N(0, \sigma^2)$

$|x_i| \leq 1, i = 1, \dots, n$; we want to choose x_1, \dots, x_n so that the least squares estimate of β has minimal variance.

Example 3 $Y_i = \beta x_i + \gamma x_i^2 + u_i, i = 1, \dots, n; u_i$ i.i.d. $N(0, \sigma^2)$;

$0 \leq x_i \leq x^*, i = 1, \dots, n$; we want to choose x_1, \dots, x_n so that the least squares estimate of $\theta = \beta x_0 + \gamma x_0^2$ has minimal variance.

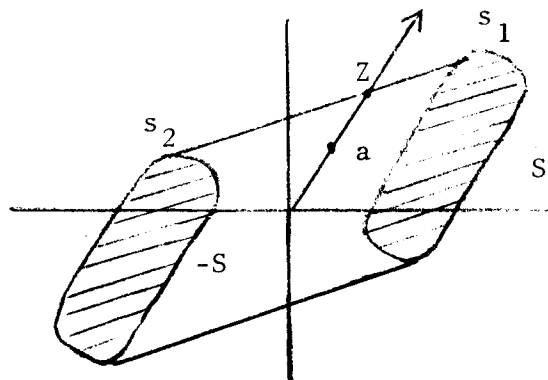
Elfving (1952, A.M.S.) derived an elegant solution to the following more general problem, which includes the two previous examples.

Example 4 $Y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + u_i, i=1, \dots, n; u_i$ i.i.d. $N(0, \sigma^2)$

$\tilde{x}_i = (x_{1i}, x_{2i}) \in$ some set S ;

We want to estimate $\theta = a_1 \beta_1 + a_2 \beta_2$.

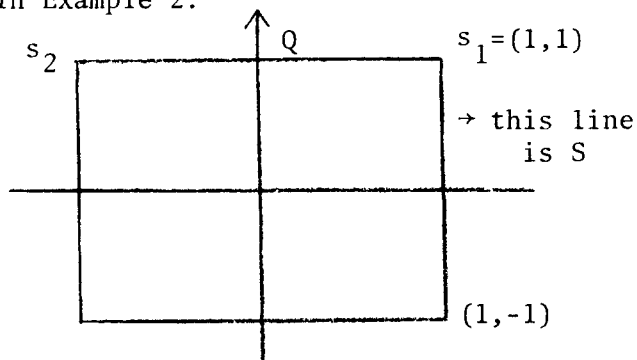
Elfving's solution: Consider the convex set generated by S and $-S$ (the pointwise reflection of S about the origin).



Draw the ray from the origin through the point $\underline{a} = (a_1, a_2)$, and consider the point \underline{z} where the ray penetrates the convex set. If $\underline{z} \in S$, repeat the level corresponding to that point n times; if $\underline{z} \in -S$, repeat the level corresponding to $-\underline{z}$ n times; if $\underline{z} \notin S, -S$, then it is a convex combination of two points $s_1, s_2 \in S \cup -S$, with weights w_1, w_2 say; repeat the levels corresponding to the points $\pm s_1, \pm s_2$ in proportions determined by these weights:

$$\text{Var}(\hat{\theta}) = \frac{\sigma^2}{n} \frac{\|\underline{a}\|^2}{\|\underline{z}\|^2}$$

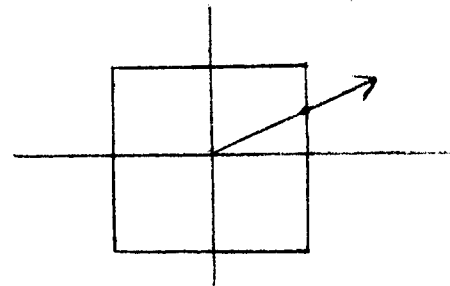
In Example 2:



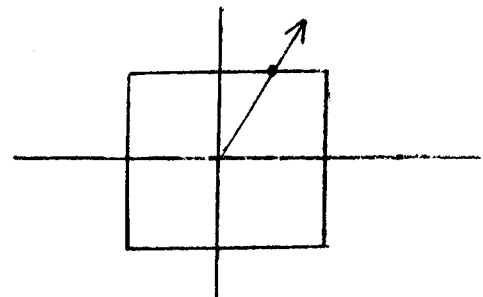
$$\underline{a} = (0,1): \underline{z} = \underline{a}$$

thus put half the observations at +1, and half at -1.

If we instead wished to estimate $\alpha + \frac{1}{2}\beta$, an optimal design puts all observations at $x = 1/2$. An alternative is to note that $(1, 1/2)$ is a weighted average of $(1,1)$ and $(1,-1)$ and to put $3/4$ of the observations at $x = 1$ and the rest at $x = -1$.



To estimate $\alpha + 2\beta$ optimally, one must put $3/4$ of the observations at $x = 1$, the rest at $x = -1$.



In Example 3, our decision depends on the particular value of x_0 : it may be optimal to place all observations at x_0 (Fig. 1), or to place some proportion of them at X and the rest at some smaller value (see Fig. 2, where the slope of the ray from 0 to Z is x_0).

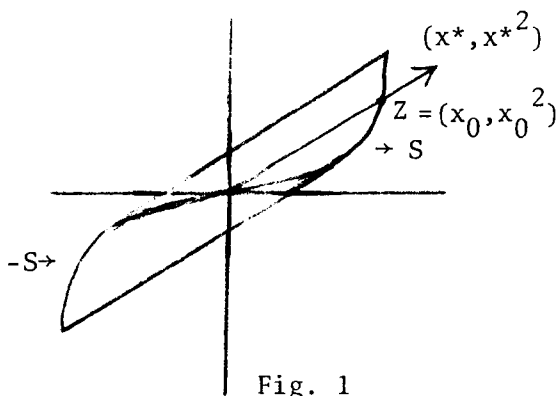


Fig. 1

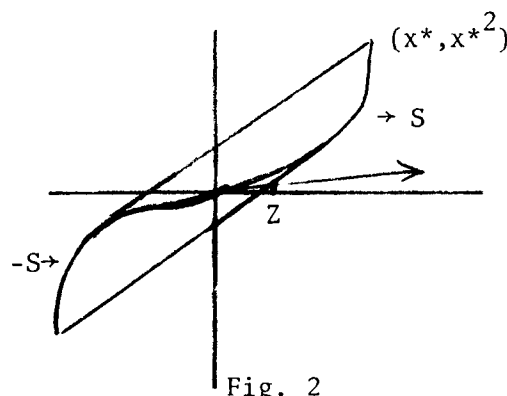


Fig. 2

Comments:

1) Suboptimal designs can be compared with optimal ones and are often preferred for practical considerations arising outside the problem statement.

2) For each choice of $\underline{x} = (x_1, x_2)^T \in S$, the corresponding information matrix is $\sigma^{-2} \underline{x} \underline{x}^T$. This has rank one which is to be expected since repetitions of this experimental level can only yield consistent estimates of one function of $\underline{\beta}$.

Non-regression problems:

Example 5 Probit (Dose response model)

We assume that the dose level x of a drug required to achieve a response is $N(\mu, \sigma^2)$, μ and σ^2 unknown. Then the probability of response to dose x is

$$p(x; \mu, \sigma^2) = \Phi\left(\frac{x-\mu}{\sigma}\right) = \Phi(w)$$

with $w = \frac{x-\mu}{\sigma}$, and where Φ is the standard normal distribution function:

$$\phi(t) = \frac{1}{\sqrt{2\pi}} e^{-t^2/2} \quad \Phi(x) = \int_{-\infty}^x \phi(t) dt$$

We want to select dose levels x_1, \dots, x_n to minimize the variance of the estimate of $\mu - 2\sigma$, the level required to achieve probability .0228 of response. The information matrix corresponding to level x becomes:

$$I_x(\theta) = \sigma^{-2} V V^T$$

where

$$\frac{\phi(w)}{\sqrt{\phi(w)(1-\phi(w))}} (1, w) = V^T$$

or

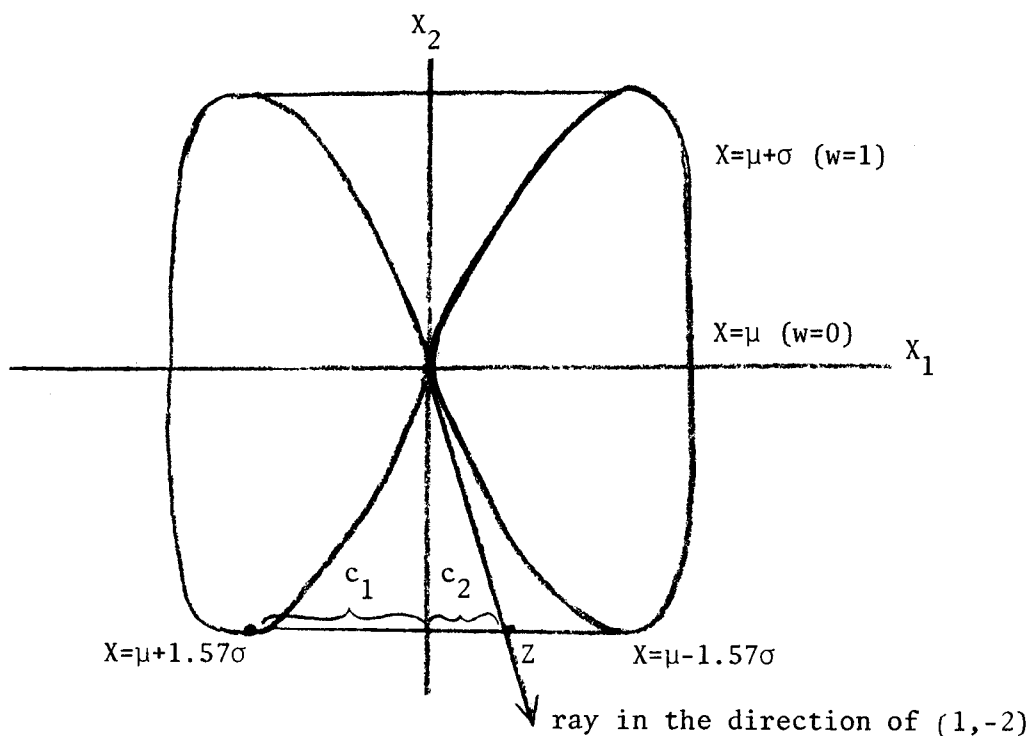
$$I_x(\theta) = \frac{\phi^2(w)}{\sigma^2 \phi(w)(1-\phi(w))} \begin{bmatrix} \bar{1} & w \\ w & w^2 \end{bmatrix}$$

$$(\phi(w) = \frac{d}{dw} \Phi(w)).$$

Since a solution to this design problem depends only on the information matrices, and the information matrices of this problem are like those of the regression problem with (X_1, X_2) replaced by

$$(\phi(w)(\phi(w)(1-\phi(w)) \sigma^2)^{-1/2}, w\phi(w)(\phi(w)(1-\phi(w)) \sigma^2)^{-1/2})$$

and (a_1, a_2) replaced by $(1, -2)$, the solution can be obtained using the Elfving method:



The solution is to put the observations at the levels $\mu \pm 1.57\sigma$ in proportions determined by c_1 and c_2 . However, the values $\mu \pm 1.57\sigma$ are unknown at the start of the experiment. We need preliminary estimates of μ and σ . Thus we must use a sequential design.

B. Optimal Design in Testing Hypotheses

In a simple hypothesis vs. simple hypothesis testing situation, we can always arrange for both error probabilities to converge to zero as $n \rightarrow \infty$, exponentially fast, say, as $e^{-n\rho}$ for some $\rho > 0$. As before, we might consider minimizing risk functions of the form

$$cn + ke^{-n\rho}$$

This is minimized by $n_0 = -\frac{1}{\rho} \log \left(\frac{c}{k\rho} \right)$. As the cost per observation c approaches zero we have $n_0 \rightarrow \infty$. Indeed $n_0 \sim -\frac{1}{\rho} \log c$, and the minimum cost is approximately $-\frac{1}{\rho} c \log c + kc$.

We can interpret $-\frac{1}{\rho} c \log c$ as the cost of sampling and kc as the cost of errors. Although these approximations are rough, they do indicate how large a

sample size to take when the cost per observation is small. Note that the cost of experimentation should be the main part of the overall risk.

A design problem to consider:

Example 6 Consider a device, the lifetime of which is exponentially distributed with failure rate θ (mean, θ^{-1}). We want to test:

$$H_1: \theta = \theta_1 \quad \text{vs} \quad H_2: \theta = \theta_2 > \theta_1$$

The following restrictions are imposed:

(i) we can only observe our sample at two time points t_1, t_2 ; thus we can determine only the number of failures in the three intervals $(0, t_1)$, (t_1, t_2) (t_2, ∞) .

(ii) letting m_i = number of components still functioning at time t_i , $i = 1, 2$; we must base our test on the statistic $T = m_1 + m_2$.

The problem is to find values t_1, t_2 which provide an optimal design for a test of the above hypotheses to be based on T , where H_1 is rejected if $T \leq k$ for some k .

In this example we may regard the statistic T as $n\bar{Y}$ where

$$\begin{aligned} Y_i &= 0 && \text{if } Z_i \leq t_1 \\ &= 1 && \text{if } t_1 < Z_i \leq t_2 \\ &= 2 && \text{if } t_2 < Z_i \end{aligned}$$

where Z_i is the lifetime of the i -th unit.

To find the asymptotic behavior of the error probabilities of this test procedure we use

Theorem 1. If $a \leq EX$, then

$$\frac{1}{n} \log\{P[\bar{X} \leq a]\} \rightarrow -\log \inf_t E[e^{t(X-a)}] = \rho(a).$$

This result may be paraphrased to state that $P[\bar{X} \leq a]$ is roughly of the order of magnitude of $e^{-n\rho(a)} = [m(a)]^n$ where

$$m(a) = \inf_t E[e^{t(X-a)}].$$

Here the term "roughly" refers to the fact that the approximation in the paraphrase may neglect factors which approach 0 or ∞ like powers of n . Indeed in many applications $P(X \leq a) \sim n^{-1/2} m^n(a)$. However, the exponential part is the most important part of this approximation for our purposes. Note that if $a \geq E(X)$, it is easy to show that $P[\bar{X} \geq a]$ is roughly of the order of $m^n(a) = e^{-n\rho(a)}$.

Our test procedure consists of rejecting H_1 if $\bar{Y} \leq a = k/n$. By selecting a between $E_{H_1}(Y)$ and $E_{H_2}(Y)$, we can show that the error probabilities of our procedure approach 0 exponentially fast in n , at rates $\rho_1(a)$ and $\rho_2(a)$, determined by a , θ_1 and θ_2 (and implicitly on t_1 and t_2). The value of a for which $\rho_1(a) = \rho_2(a) = \rho_0$ will yield error probabilities $\alpha = \epsilon_1$ and $\beta = \epsilon_2$ for which $\alpha + \lambda\beta$ is approximately minimized and for which roughly $\epsilon_1 \sim \epsilon_2 \sim e^{-n\rho_0}$. Our optimal design is obtained by selecting t_1 and t_2 appropriately to maximize ρ_0 .

Suppose now that we decide to remove the restriction on the use of T . Given the design, i.e. the choice of t_1 and t_2 , we may wish to use the efficient likelihood-ratio test. However, the likelihood-ratio test is also a test of the form: reject H_0 if $\bar{Y} \leq a$.

To see this we note that the likelihood-ratio test for testing $H_1: f(x) = f_1(x)$ vs $H_2: f(x) = f_2(x)$ consists of rejecting H_1 if

$$\lambda = \prod_{i=1}^n \frac{f_1(X_i)}{f_2(X_i)} \leq k_1$$

or if

$$\sum Y_i \leq k_2 = \log k_1$$

or if

$$\bar{Y} \leq a = k_2/n$$

where

$$Y_i = \log[f_1(X_i)/f_2(X_i)] .$$

In the Bayesian context where the wrong decision under H_i costs r_i and H_i has prior probability π_i , $k_1 = \pi_2 r_2 / \pi_1 r_1$ and $a \rightarrow 0$ as $n \rightarrow \infty$. Then it is easy to see that for the Bayes tests the error probabilities approach zero roughly like

$$m_0^n = e^{-n\rho_0}$$

where

$$m_0 = \inf_{0 \leq t \leq 1} \int f_1^t(x) f_2^{1-t}(x) dx$$

if f_1 and f_2 correspond to continuous distributions. The integrals are replaced by a sum if the distributions are discrete.

Thus each design choice in our less restricted problem consists of selecting t_1 and t_2 appropriately. That is, for each t_1, t_2 , there is a corresponding f_1 and f_2 and m_0 , both of which are discrete distributions with three possible values. We select t_1 and t_2 to minimize m_0 or maximize $\rho_0 = -\log m_0$.

We shall not pursue further the numerical details for this problem.

III. Sequential Design of Experiments

A. Sequential Analysis for Testing Simple Hypotheses versus Simple Alternatives

Wald (1947) introduced the SPRT (Sequential Probability Ratio Test) for testing $H_1: f(x) = f_1(x)$ versus $H_2: f(x) = f_2(x)$ where $f(x)$ is the distribution function of the i.i.d. random variables X_1, X_2, \dots . This test procedure is to

Reject H_1 if $\lambda_n \leq B$

Accept H_1 if $\lambda_n \geq A$

and continue sampling if $B < \lambda_n < A$

where

$$\lambda_n = \prod_{i=1}^n \frac{f_1(X_i)}{f_2(X_i)}$$

is the sequential probability-ratio and $B < 1 < A$. Then

$$S_n = \log \lambda_n = \sum_{i=1}^n Y_i$$

where the

$$Y_i = \log \frac{f_1(X_i)}{f_2(X_i)}$$

are i.i.d. observations on Y .

The SPRT has several basic properties:

1) The SPRT leads to termination with probability one. More precisely, let N be the sample size n when λ_n first fails to be within (B, A) . We let $N = \infty$ if $B < \lambda_n < A$ for all n . As long as $P_{H_1}(Y=0) < 1$ for $i = 1, 2$, i.e. $f_1(x)$ and $f_2(x)$ represent distinct probability distributions, $P_{H_i}(N=\infty) = 0$, $i=1, 2$.

2) The error probabilities $\alpha = \varepsilon_1 = P_{H_1}(\text{Reject } H_1)$ and $\beta = \varepsilon_2 = P_{H_2}(\text{Accept } H_1)$

satisfy

$$\frac{\alpha}{1-\beta} \leq B$$

and

$$\frac{1-\alpha}{\beta} \geq A .$$

Indeed, these inequalities are approximations under conditions where the "excess" is small. The excess corresponds to the overshoot beyond A or B of the probability ratio at termination of sampling. These approximations, for which bounds can be constructed, state that

$$\alpha \approx \frac{B(A-1)}{A-B}$$

$$\beta \approx \frac{1-B}{A-B}$$

3) The expected sample size $E_{H_1}(N)$ satisfies

$$E_{H_1}[S_N] = E_{H_1}(N)E_{H_1}(Y) .$$

Assuming that the excess is negligible, $S_N \approx \log A$ when H_1 is accepted and

$S_N \approx \log B$ when H_1 is rejected and hence

$$E_{H_1}(S_N) \approx \alpha \log B + (1-\alpha) \log A$$

$$E_{H_2}(S_N) \approx (1-\beta) \log B + \beta \log A .$$

Moreover,

$$E_{H_1}(Y) = E_{H_1} \left\{ \log \left[\frac{f_1(X)}{f_2(X)} \right] \right\} = I(f_1, f_2)$$

and

$$-E_{H_2}(Y) = E_{H_2} \left\{ \log \left[\frac{f_2(X)}{f_1(X)} \right] \right\} = I(f_2, f_1)$$

are the Kullback-Leibler information numbers for discriminating between f_1 and f_2 .

Thus

$$E_{H_1}(N) \approx \frac{1}{I(f_1, f_2)} \{ \alpha \log B + (1-\alpha) \log A \}$$

and

$$E_{H_2}(N) \approx \frac{1}{I(f_2, f_1)} \{ (\beta-1) \log B - \beta \log A \}$$

4) The sequential probability ratio test is optimal.

Given a SPRT with error probabilities α and β , for any other test with error probabilities $\alpha' \leq \alpha$ and $\beta' \leq \beta$ both expected sample sizes under H_1 and H_2 are at least equal to those for the SPRT. Moreover, the Bayes sequential tests are SPRT's if the cost of sampling is linear in sample size.

B. Asymptotic Behavior of the Bayes SPRT

Suppose that the cost of the wrong decision when H_i is true is r_i and that the cost per observation is c . Then the Bayes sequential test as $c \rightarrow 0$ involves taking a large number of observations. This means that $A \rightarrow \infty$ and $B \rightarrow 0$. Assuming that $\log A$ and $-\log B$ are of the same order of magnitude, one may approximate

$$\alpha \approx B, \beta \approx A^{-1}, E_{H_1}(N) \approx \frac{\log A}{I(f_1, f_2)}, \text{ and } E_{H_2}(N) \approx \frac{-\log B}{I(f_2, f_1)}.$$

Then minimizing the Bayes Risk

$$R = \pi_1 R_1 + \pi_2 R_2$$

where the risks

$$R_1 = cE_{H_1}(N) + r_1 \alpha$$

$$R_2 = cE_{H_2}(N) + r_2 \beta$$

and π_i are the prior probabilities of the H_i , leads to

$$A \approx \frac{\pi_2}{\pi_1} \frac{r_2 I(f_1, f_2)}{c}, \quad B \approx \frac{\pi_2}{\pi_1} \frac{c}{r_1 I(f_2, f_1)},$$

$$\alpha \sim c, \quad \beta \sim c,$$

$$E_{H_1}(N) \approx \frac{-\log c}{I(f_1, f_2)}, \quad E_{H_2}(N) \approx \frac{-\log c}{I(f_2, f_1)},$$

$$R_1 \approx \frac{-c \log c}{I(f_1, f_2)}, \quad \text{and} \quad R_2 \approx \frac{-c \log c}{I(f_2, f_1)}.$$

Thus the main part of the expected cost for the Bayes Procedure comes from the cost of sampling and is of the order of magnitude of $-c \log c$. Furthermore, this cost is inversely proportional to the Kullback-Leibler Information Number and relatively insensitive to the prior probabilities and the costs of incorrect decisions.

To see the relevance of this result for the problem of design of experiments, suppose that there are available two experiments for sequentially testing H_1 vs H_2 . Under the first experiment we observe X_1, X_2, \dots with information numbers $I(f_1, f_2) = I_1$, and $I(f_2, f_1) = I_2$. Suppose that for the second experiment we observe X_1^*, X_2^*, \dots with information numbers I_1^* and I_2^* .

If $I_1 > I_1^*$ and $I_2 > I_2^*$, (and the same cost c per observation) it seems clear that we should prefer to observe the X_i . If $I_1 > I_1^*$ but $I_2 < I_2^*$, the situation is not as clear. Then X_1, X_2, \dots are to be preferred if H_1 is true but X_1^*, X_2^*, \dots if H_2 is true. If we knew which hypothesis were true, there would be no need to experiment. However, as data cumulate, one may be still unsure enough about which hypothesis is true to invest in further observations without being sure enough to stop and make a final decision. In that event it makes sense to select X if $\lambda_n > 1$ and X^* if $\lambda_n < 1$.

This strategy would involve shifting from X to X^* or vice versa as λ_n changes during experimentation.

Another insight to be gained from the above asymptotic analysis is the following. Let π_{1n} and π_{2n} be the posterior probabilities after n observations X_1, X_2, \dots, X_n . Then Bayes Theorem states that

$$\frac{\pi_{1n}}{\pi_{2n}} = \frac{\pi_1}{\pi_2} \lambda_n$$

Thus the asymptotic Bayes sequential rule corresponds to stopping and selecting H_1 when π_{2n} goes below some number roughly of the order of magnitude of c . Similarly we stop and select H_2 when π_{1n} goes below some such number.

C. A Simple Composite Hypothesis Testing Problem.

The simplest composite problem will shed some light on our problem. Let us test $H_1: \theta = \theta_1$ vs $H_2: \theta = \theta_2$ or $\theta = \theta_3$. Let $\underline{\pi} = (\pi_1, \pi_2, \pi_3)$ represent the prior probability distribution.

In the Bayesian framework we may compute the posterior distribution

$\underline{\pi}_n = (\pi_{1n}, \pi_{2n}, \pi_{3n})$ where

$$\frac{\pi_{2n}}{\pi_{1n}} = \frac{\pi_2}{\pi_1} e^{-S_{12n}} \quad \frac{\pi_{3n}}{\pi_{1n}} = \frac{\pi_3}{\pi_1} e^{-S_{13n}}$$

$$S_{12n} = \sum_{i=1}^n Y_{12i} = \sum \log \frac{f(X_i | \theta_1)}{f(X_i | \theta_2)}$$

$$S_{13n} = \sum_{i=1}^n Y_{13i} = \sum \log \frac{f(X_i | \theta_1)}{f(X_i | \theta_3)}$$

Suppose now that $H_1: \theta = \theta_1$, is true. Then

$$E_{\theta_1}(Y_{12i}) = I(\theta_1, \theta_2) > 0$$

$$E_{\theta_1}(Y_{13i}) = I(\theta_1, \theta_3) > 0$$

and, assuming $\pi_1 > 0$, π_{2n} and π_{3n} will approach zero roughly as $\exp[-nI(\theta_1, \theta_2)]$ and $\exp[-nI(\theta_1, \theta_3)]$. The posterior probability of $H_2 = \pi_{2n} + \pi_{3n}$ will approach zero as the sum of these, i.e., like $\exp[-nI(\theta_1)]$ where

$$I(\theta_1) = \min[I(\theta_1, \theta_2), I(\theta_1, \theta_3)]$$

When H_1 is true the experimenter should select an experiment which maximizes $I(\theta_1)$. The following example is informative.

Example 7 We have three experiments with $I(\theta_1, \theta_i)$ described in the following table.

	e_1	e_2	e_3
$I(\theta_1, \theta_2)$	6	3	2
$I(\theta_1, \theta_3)$	2	3	6
$I(\theta_1)$	2	3	2

Of these three experiments e_2 will maximize $I(\theta_1)$. However, if the statistician were to randomize and select e_1 and e_3 each with probability 1/2, we would have $I(\theta_1, \theta_2) = I(\theta_1, \theta_3) = I(\theta_1) = 4$. This suggests a potential advantage of selecting experiments by using randomization.

D. A Proposed Plan for the Sequential Design of Experiments

Assembling the suggestions of the preceding section we propose a procedure for the problem of sequential design of experiments formulated as follows:
Let E be a set of available elementary experiments e whose outcome X_e has

density $f(x|\theta, e)$ with respect to a measure μ_e . Let $H_1: \theta \in \Theta_1$, and $H_2: \theta \in \Theta_2$ be two hypotheses where Θ_1 and Θ_2 are disjoint sets whose union Θ represents the set of possible states of nature. Let $r(\theta)$ be the cost of deciding incorrectly and let c be the cost per observation. After each observation a decision is made to stop or continue sampling. If sampling is stopped, either H_1 or H_2 is accepted. If sampling is continued a new experiment $e_{n+1} \in E$ is selected. While its choice may depend on the past data (X_1, X_2, \dots, X_n) once the experiment is selected its outcome is independent of the past.

To describe the proposed solution, let E^* be the set of randomized experiments generated by E . For $e \in E^*$,

$$Y(\theta, \phi, e) = \log[f(X_e|\theta, e)/f(X_e|\phi, e)]$$

has expectation $I(\theta, \phi, e)$ and variance $v(\theta, \phi, e)$ when θ is the true state of nature. Let $Y_n(\theta, \phi) = \log[f(X_n|\theta, e_n)/f(X_n|\phi, e_n)]$ and

$$S_n(\theta, \phi) = \sum_{i=1}^n Y_i(\theta, \phi)$$

and $\hat{\theta}_n$ be the maximum-likelihood estimate (m.l.e.) of θ based on X_1, X_2, \dots, X_n . Let $h(\theta) = \Theta_i$ and $a(\theta) = \Theta - \Theta_i$ if $\theta \in \Theta_i$, $i=1,2$ and let $\tilde{\theta}_n$ be the m.l.e. of θ based on X_1, \dots, X_n when θ is restricted to $a(\hat{\theta}_n)$. Here $h(\theta)$ is identified with the "hypothesis of θ " and $a(\theta)$ with the "alternative to the hypothesis of θ ", and $\tilde{\theta}_n$ is the m.l.e. on the alternative to the hypothesis of $\hat{\theta}_n$.

We describe our sequential procedure, Procedure A: stop sampling and accept the hypothesis $h(\hat{\theta}_n)$ if

$$S_n(\hat{\theta}_n, \tilde{\theta}_n) > -\log c.$$

Otherwise select $e_{n+1} \in E^*$ to maximize

$$\inf_{\phi \in a(\hat{\theta}_n)} I(\hat{\theta}_n, \phi, e).$$

On the basis of our motivation one would anticipate that the procedure would be asymptotically optimal and that its risk would be asymptotically

$$R(\theta) \approx -c \log c / I(\theta)$$

where

$$I(\theta) = \sup_{e \in E^*} \inf_{\phi \in a(\theta)} I(\theta, \phi, e).$$

Indeed, the following theorem holds:

Theorem If E and Θ are finite and for every $e \in E$, $I(\theta, \phi, e) > 0$ and $v(\theta, \phi, e) < \infty$ for all θ and $\phi \in \Theta$ with $\theta \neq \phi$, and $I(\theta) > 0$ then the risk $R(\theta)$ for procedure A satisfies

$$R(\theta) \leq [1 + o(1)] \left\{ \frac{-c \log c}{I(\theta)} \right\} \text{ for all } \theta$$

and any procedures for which the risk $R^*(\theta) = O(-c \log c)$ for all θ satisfies

$$R(\theta) \geq [1 + o(1)] \left[\frac{-c \log c}{I(\theta)} \right] \text{ for all } \theta.$$

In short, this theorem states that procedure A is optimal in the sense that to do substantially better for any state of nature, one must do worse by an order of magnitude for some states of nature.

This result has been generalized considerably. It has two major shortcomings. It breaks down when points of θ_1 and θ_2 can get close in the sense of $I(\theta_1, \theta_2) \rightarrow 0$. It is not very dependable in the early stages of experimenting and requires supplementation for problems which will require moderate sample sizes.

IV. Continuous Time Sequential Problems

A. Sequential Problems

We list three problems.

1) Deciding sign of the mean of a normal distribution.

Let $X_1, X_2, \dots, X_n \dots$ be i.i.d. $N(\mu, \sigma^2)$ with σ^2 known. It is desired to test $H_1: \mu > 0$ vs $H_2: \mu \leq 0$. The SPRT was designed to test a simple hypothesis versus a simple alternative and could be modified for composite hypotheses. However, such modifications are inadequate for this case where points of the alternative hypotheses are close to one another. The classical approach consists of assuming the existence of an indifference zone (μ_1, μ_2) within which it doesn't matter what decision is made and testing $H_1^*: \mu = \mu_1$ vs $H_2^*: \mu = \mu_2$. This approach leads to unduly large sample sizes for $\mu = (\mu_1 + \mu_2)/2$ and is theoretically inadequate if the cost of sampling is small and there is a positive cost of wrong decision for all $\mu \neq 0$.

We complete the statement of the problem in a decision theoretic form and apply a Bayesian approach to state a clean cut optimization problem. The cost of sampling is c per observation. The cost of a wrong decision is $r(\mu) = k(\mu)$. Finally we assume that μ has the $N(\mu_0, \sigma_0^2)$ prior distribution. We seek a sequential decision procedure which minimizes the expected cost.

2) One armed bandit.

Let X_1, X_2, \dots, X_M be i.i.d. $N(\mu, \sigma^2)$, with σ^2 known. A player receives $X_1 + X_2 + \dots + X_n$ if he stops at $N = n \leq M$, where the stopping time N may be based on the past data $\underline{X}_N = (X_1, \dots, X_N)$. The unknown mean μ has the $N(\mu_0, \sigma_0^2)$ prior distribution.

This problem is a normal version of the one armed bandit problem where a

player may play up to M times on a one armed bandit in an environment where most "bandits" may be unfavorable but it is known that some are favorable.

3) Canonical problem.

Let Y_n be a process starting at y_0 for $n = n_0$, where n_0 is a negative integer. As n increases $Y_{n+1} = Y_n + u_n$ where the u_n are independent $N(0,1)$ random variables. A player may stop with no payoff or at a cost of one, increase n by one. If he continues till $n = 0$, the game ends and he receives 0 if $Y_0 \geq 0$ and Y_0^2 if $Y_0 \leq 0$.

B. Bayesian Analysis

In problems 1 and 2, the data $X_n = (X_1, X_2, \dots, X_n)$ and the prior $N(\mu_0, \sigma_0^2)$ give rise to the posterior distribution.

$$L(\mu | X_n) = N(Y_n, s_n)$$

where

$$Y_n = \frac{\mu_0 \sigma_0^{-2} + (X_1 + \dots + X_n) \sigma^{-2}}{\sigma_0^{-2} + n \sigma^{-2}}$$

and

$$s_n = (\sigma_0^{-2} + n \sigma^{-2})^{-1} \quad \text{is decreasing in } n.$$

It is possible to show that as data cumulate, Y_n , the Bayesian estimate of μ , changes according to independent increments, and

$$L(Y_n - Y_m | Y_m) = N(0, s_m - s_n) \quad \text{for } n \geq m \geq 0.$$

In problem 1, the posterior risk upon stopping at $N = n$ is

$$cn + E\{k | \mu | X_n\} = d_1(Y_n, s_n)$$

where

$$d_1(y, s) = ks^{1/2} \psi(ys^{-1/2}) + c\sigma^2 s^{-1} - c\sigma_0^{-2} \sigma^2$$

and

$$\psi(\mu) = \begin{cases} \phi(u) = u[1 - \Phi(u)] & u \geq 0 \\ \phi(u) + u\Phi(u) & u \leq 0. \end{cases}$$

Thus our sequential Problem 1 has been reduced to the stopping problem of observing the Gaussian stochastic process $\{Y_n; n = 1, 2, \dots\}$ of independent increments and stopping with cost $d_1(y, s)$ if stopping takes place at $Y_n = y, s_n = s$. The object is to find the stopping rule which minimizes $E\{d_1[Y_n, s_n]\}$.

A similar analysis shows that Problem 2 also reduces to a stopping problem with stopping cost

$$d_2(y, s) = \sigma^2 [\mu_0 \sigma_0^{-2} - y/s] \quad \text{for } 0 \leq n \leq M.$$

Here stopping is enforced at M if it has not taken place earlier. Problem 3 is also easily seen to be a stopping problem with $\{Y_n; n = n_0, n_0 + 1, \dots, 0\}$, $s = n$, stopping enforced at $n = 0$, and stopping cost

$$d_3(y, s) = \begin{cases} -n_0 - s - y^2 & \text{if } s = 0 \text{ and } y < 0 \\ -n_0 - s & \text{otherwise} \end{cases}$$

C. Continuous Time Version of Stopping Problems

Each of the three problems of Section A has a continuous time version. These versions are of interest in themselves. The solutions of the discrete time problems are approximated by those of the continuous versions which are more readily subject to mathematical analysis.

Problem 1*. We observe $X(t)$, a Wiener process (Gaussian process of independent increments) with drift μ per unit time and variance σ^2 per unit time; i.e., $E[dX(t)] = \mu dt$, $\text{Var}[dX(t)] = \sigma^2 dt$. Once again the prior distribution is $N(\mu_0, \sigma_0^2)$ and the cost of wrong decision is $r(\mu) = k(\mu)$.

Here the posterior distribution of μ given $X(t')$, $0 \leq t' \leq t$ is

$$L(\mu|X(t'), 0 \leq t' \leq t) = N(Y, s)$$

where

$$Y = Y(s) = \frac{\mu_0 \sigma_0^{-2} + X(t) \sigma^{-2}}{\sigma_0^{-2} + t \sigma^{-2}}$$

and

$$s = (\sigma_0^{-2} + t \sigma^{-2})^{-1}.$$

The process $Y(s)$ is a Wiener process in the $-s$ scale starting at

$$(Y(s), s) = (\mu_0, \sigma_0^2), \text{ i.e., } E[dY(s)] = 0, \text{ and } \text{Var}[dY(s)] = -ds.$$

The Bayes strategy for this problem is the solution of the continuous time stopping problem involving $\{Y(s): \sigma_0^2 \geq s > 0\}$ with stopping cost $d_1(y, s)$.

The transformation $s^* = a^2 s$, $Y^* = aY$, permits one to normalize this problem to one with stopping cost

$$d_1^*(y, s) = s^{-1} + s^{1/2} \psi(y, s^{-1/2}).$$

Problem 2*. The continuous time version of the one armed bandit problem is similar to Problem 1* except for the stopping cost $d_2(y, s)$ which can be normalized to

$$d_2^*(y, s) = -y/s \quad \text{for } s \geq 1$$

with stopping imposed at $s = 1$.

Problem 3*. The continuous time version of this problem involves observing $\{Y(s): s_0 \geq s \geq 0\}$ with $Y(s_0) = y_0$, stopping imposed at $s = 0$ and stopping cost

$$d_3^*(y, s) = \begin{cases} -s - y^2 & \text{if } s = 0 \text{ and } y < 0 \\ -s & \text{otherwise} \end{cases}$$

D. Free Boundary Problem

In this section we shall relate the optimal solution of continuous time stopping problems to that of a free boundary problem involving the heat equation.

First let $b(y,s)$ represent the risk or expected cost for a specified stopping rule given that the stochastic process Y has reached $Y(s) = y$. That is

$$b(y,s) = E\{d[y(S),S] | Y(s) = y\}$$

where $S > s$ is a stopping time (identified with the rule). Let $\rho(y,s)$ be the optimal risk achievable given $Y(s')$ for $s' \geq s$ and $Y(s) = y$. Since Y is a process of independent increments, ρ depends on $Y(s')$ only through (y,s) . Then $\rho(y,s) \geq d(y,s)$ and for a procedure to be optimal it must call for stopping only when $\rho(y,s) = d(y,s)$. Thus we may confine attention to stopping rules which are determined by a stopping set S and the complementary continuation set C in the (y,s) plane.

Clearly, $b(y,s) = d(y,s)$ for $(y,s) \in S$. For (y,s) an inner equation point of C ,

$$b(y,s+\delta) = E\{b(Y(s),s) | Y(s+\delta) = y\}$$

together with the fact that the conditional distribution of $Y(s)$ is Normal with mean y and variance δ yields the heat equation

$$\frac{1}{2} b_{yy} = b_s \quad \text{for } (y,s) \in C.$$

A somewhat more delicate analysis shows that the optimality condition for the stopping rule reduces to smoothness of b across the boundary of C . Thus the problem of finding an optimal stopping set S is related to that of finding the solution (ρ,C) of the following free boundary problem (f.b.p.)

$$\frac{1}{2} \rho_{yy}(y, s) = \rho_s(y, s) \text{ on } C$$

$$\rho(y, s) = d(y, s) \text{ on } S$$

$$\rho_y(y, s) = d_y(y, s) \text{ on } B$$

where B is the boundary of C and S is the complement of C .

The relation between the optimization problem and the free boundary problem permits one to apply the methods of partial differential equations. Three major techniques have been found to be useful, when closed form solutions are not available.

One of these due to Bather (1962) is to generate solutions of the heat equation and to find what stopping problems they solve. By relating these stopping problems to ours, we find bounds on the solution to our problem.

A simpler variation is the following:

Let us consider the normalized version of Problem 1* with

$$d_1^*(y, s) = s^{-1} + s^{1/2} \psi(ys^{-1/2})$$

A special solution of the heat equation is

$$u_1(y, s) = Ks^{-1/2} \phi(ys^{-1/2}) .$$

For $K > 3(\pi/2)^{1/3}$, the equation

$$d_1^*(y, s) = u_1(y, s)$$

determines the curves $\pm y_1^*(s)$ for $s \geq s_k$ and $y_1^*(s_k) = 0$. Then let

$$b_1(y, s) = u_1(y, s) \quad \text{on } C = \{(y, s) : s > s_k, |y| < y_1^*(s)\}$$

$$b_1(y, s) = d_1^*(y, s) \quad \text{on } S$$

and b is the risk function corresponding to the stopping rule determined by the continuation set C . This implies that the optimal risk $\rho_1^* \leq b_1$. Thus any point (y,s) for which $b_1(y,s) < d_1^*(y,s)$ is one for which $\rho_1^*(y,s) < d_1^*(y,s)$ and hence a point of the optimal continuation set.

The function $b_2(y,s)$ defined by

$$b_2(y,s) = u_2(y,s) = -y \quad \text{for } s > 1, y > 0 \quad (C_2^{**})$$

$$= d_2^*(y,s) = -y/s \quad s \geq 1, y \leq 0 \quad (S_2^{**})$$

is a solution of the ordinary boundary value problem for $d_2^*(y,s)$. Hence b_2 is the risk for the suboptimal procedure C_2^{**} . For any point (y,s) for which $b_2 < d_2^*$ the optimal risk ρ_2 for Problem 2* will certainly satisfy $\rho_2 < d_2^*$ and (y,s) will be in C_2^* , the optimal continuation set for d_2^* . But $b_2 < d_2^*$ for $y > 0$ and $s > 1$. Hence

$$C_2^* \supset \{(y,s): s > 1, y > 0\}.$$

A second method, related to the first, consists of developing asymptotic expansions for large s and for small s . Thus for Problem 1*, the symmetric solution has the property

$$\tilde{y}_1(s) s^{-1/2} \sim \{\log s^3 - \log 8\pi - 6(\log s^3)^{-1} + \dots\}^{1/2} \quad \text{as } s \rightarrow \infty$$

$$\tilde{y}_1(s) s^{-1/2} \sim (0.25) s^{3/2} \{1 - (1/12) s^3 + (7/15 \cdot 16) s^6 - \dots\} \quad \text{as } s \rightarrow 0$$

where $\pm \tilde{y}_1(s)$ are the boundaries of the optimal continuation set. For Problem 2*, the optimal continuation set lies above the curve $\tilde{y}_2(s)$ for which

$$\begin{aligned} \tilde{\beta} = \Phi[s^{-1/2} \tilde{y}_2(s)] &\approx 2s^{-1} && \text{for } s \rightarrow \infty \\ &\approx 1/2 && \text{for } s \rightarrow 1. \end{aligned}$$

Note that $\tilde{\beta}$ can be interpreted as a nominal significance level which leads to termination and s^{-1} represents that proportion of the total potential information that is available at the time.

$$\rho(y,s) = u_3(y,s) \equiv y^2 - s \quad \text{for } y \leq 0$$

$$\rho(y,s) = -s \quad \text{for } y \geq 0.$$

Then ρ is a solution of the f.b.p. and the optimal procedure is to stop whenever $Y(s) \geq 0$.

A third approach is to use numerical methods. These will be described in the next section.

In the meantime let us observe that Problem 3* has a trivial solution which is given by

E. Numerical Approximation

The continuous time stopping problems were introduced to furnish analytical descriptions of approximations to the discrete time problem. Closed form solutions are rarely available and numerical methods are desirable. One natural numerical method consists of approximating the continuous time problem by a related discrete time version where stopping is permitted only at closely spaced discrete time points. The latter problem may be solved using backward induction.

We seem to have gone in a logical circle. This is not completely the case. First, the excursion to continuous time has permitted us to construct bounds and expansions. Second, a single discrete approximation can be used to approximate the continuous problem which in turn is an approximation to a whole class of discrete time problems with a given normalized version.

A direct approach to this numerical solution by backward induction has two shortcomings. First, this approach involves many numerical integrations

corresponding to expectations with respect to normal densities. Second, a doubling of accuracy involves 2^3 times as many calculations.

We shall reduce the integration problem by replacing the Gaussian process by the discrete

$$Y(s-\delta) = Y(s) + \sqrt{\delta} \quad \text{w.p. } 1/2$$

$$Y(s) - \sqrt{\delta} \quad \text{w.p. } 1/2.$$

This process has independent increments with mean 0 and variance 1 per unit time. The backward induction equation which replaces the numerical integrations when the Wiener process is used is the relatively simple

$$\rho(y, s+\delta) = \min[d(y, s+\delta), \frac{\rho(y+\sqrt{\delta}, s) + \rho(y-\sqrt{\delta}, s)}{2}] .$$

When δ is small, this process will approximate the continuous time solution.

The second problem is alleviated considerably by two facts. For all the problems with which I have dealt, this backward induction gives surprisingly good approximations with very coarse grids $(\delta, \sqrt{\delta})$ in the s and y scales. Secondly there is a correction for the discreteness which effectively increases the accuracy of the method considerably, thereby replacing the cube factor for refinement by something very close to linear. Thus this method is very useful for practical purposes although it still leaves something to be desired for those who require highly accurate solutions.

The correction factor can be derived by an argument which points out that to a "small statistician" located near the boundary of the continuous time solution, the stopping problem very much resembles Problem 3*. Then the relation between the corresponding discrete and continuous time solutions of Problem 3* provides a good approximation to the difference between the discrete and continuous time solutions of the general stopping problem.

We pose Problem 3**. Let $Y^{**}(s)$ be such that

$$Y^{**}(s) = Y^{**}(s+1) \pm 1 \quad \text{w.p. } 1/2, \quad s = 0, 1, 2, \dots$$

and let $d_3^*(y,s)$ be the stopping cost. Find the optimal stopping procedure and the associated risk. The optimal strategy is to stop when $Y^{**}(s) \geq \tilde{y}_3^{**}(s)$ where $\tilde{y}_3^{**}(s)$ is monotone decreasing in s and approaches $-1/2$.

Hence a first correction to the solution of a continuous time solution of a stopping problem consists of correcting the boundary of the discrete approximation by expanding the continuation region boundary by $0.5\sqrt{\delta}$ along the y direction.

This correction is rather crude because it represents $1/2$ the grid size in the y direction, and there is a little vagueness about where the boundary belongs if two neighboring grid points are such that one is a stopping point and the other a continuation point for the discrete problem. This difficulty is resolved by studying the nearby risk values and comparing the optimal continuous time and discrete time risks for the canonical problem as $s \rightarrow \infty$. The continuous time optimal risk is $-y^2 - s$ for $y \leq 0$ and $-s$ for $y \geq 0$. The discrete time optimal risk, for large s , behaves like $v(y) - s$ for $y \leq -.5$ and $-s$ for $y \geq -.5$ where

$$v(y) = -y^2 + \inf\{(y+j)^2 : j \text{ an integer}\}.$$

This leads to the approximation

$$\tilde{y}(s_i) = \tilde{y}^{**}(s_i) - (1-u)\delta^{1/2}$$

where

$$u = (D_1 - 4D_0)/2(D_1 - D_0)$$

$$D_0 = \rho^{**}(y_0, s) - d^*(y_0, s)$$

$$D_1 = \rho^{**}(y_1, s) - d^*(y_1, s).$$

and y_0 and y_1 are the two continuation points closest to the optimal boundary of the discrete problem.

REFERENCES

- Bather, J. A., Bayes procedure for deciding the sign of a normal mean, *Proc. Cambridge Philos. Soc.* 58 (1962), pp 599-620.
- Breakwell, J. V. and H. Chernoff, Sequential tests for the mean of a normal distribution II (large t), *Ann. Math. Statist.* 35 (1963), pp 162-173.
- Chernoff, H., Locally optimal designs for estimating parameters, *Ibid.*, 24 (1943), pp 586-602.
- _____, A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations, *Ibid.*, 23 (1952), pp 493-507.
- _____, Sequential analysis and optimal design, *SIAM*, Philadelphia, Pennsylvania, 1972.
- _____, Sequential design of experiments, *Ann. Math. Statist.* 30 (1959), pp 755-770.
- _____, Sequential tests for the mean of a normal distribution, *Proc. Fourth Berkeley Symp. Math. Statist.* 1 (1961), pp 79-91.
- _____, Sequential tests for the mean of a normal distribution III (small t), *Ann. Math. Statist.* 36 (1965), pp 28-54.
- _____, Sequential tests for the mean of a normal distribution IV (discrete case), *Ibid.*, 36 (1965), pp 55-68.
- _____, and A. J. Petkau, An optimal stopping problem for sums of dichotomous random variables, *Ann. Probabil.* 4 (1976), pp 875-889.
- _____, and S. N. Ray, A Bayes sequential sampling inspection plan, *Ibid.*, 36 (1965), pp 1387-1407.
- Elfving, G., Optimum allocation in linear regression theory, *Ann. Math. Statist.* 23 (1952), pp 255-262.
- Wald, A., *Sequential Analysis*, John Wiley, New York, 1947.