

TO APPEAR IN  
Advances in Applied Probability  
(1980) 12 No. 3.

**TWO-ARMED BANDITS WITH A GOAL;**

**I: ONE ARM KNOWN**

by

**Donald A. Berry\***

and

**Bert Fristedt\*\***

**University of Minnesota**

**Technical Report No. 344**

August, 1979

Corrections: October 1979

\* This author's research sponsored by the NSF under Grant No. 78-02694.

\*\* This author's research sponsored by the NSF under Grant No. 74-05786 A02.

Two-armed bandits with a goal;

I: One arm known

by

Donald A. Berry and Bert Fristedt

ABSTRACT

One of two random variables,  $X$  and  $Y$ , can be selected at each of a possibly infinite number of stages. Depending on the outcome, one's fortune is either increased or decreased by one. The probability of increase may not be known for either  $X$  or  $Y$ . The objective is to increase one's fortune to  $G$  before it decreases to  $g$ , for some integral  $g$  and  $G$ ; either may be infinite.

In the current part of the paper, the distribution of  $X$  is unknown and that of  $Y$  is known. We characterize the situations in which optimal strategies exist and, for certain kinds of information concerning  $X$  and  $Y$ , we characterize optimal sequential strategies for choosing to observe  $X$  and  $Y$ .

In Part II (Berry and Fristedt 1980), it is known that either  $X$  or  $Y$  has probability  $\alpha$  of increasing the current fortune by one and the other has probability  $\beta$  of increasing the fortune by one, where  $\alpha$  and  $\beta$  are known, but which goes with  $X$  is not known.

Key words and phrases: Achieving a goal, two-armed bandits, how to gamble if you must, gambler's ruin, sequential decisions, Bayesian decision making, sequential medical treatments, stochastic control, optimal dynamic designs.

Acknowledgement: This problem and its applicability to certain medical treatment problems were suggested to us by Professor William D. Sudderth.

1. Introduction. A gambler has an integral fortune  $k$  and is allowed to make a sequence of unit bets. He is to continue betting until his fortune becomes either of the integral goals  $G$  or  $g$ , with  $g < k < G$ . Each bet is made by choosing to observe one of two random variables,  $X$  and  $Y$ , where

$$P(X = 1) = 1 - P(X = -1) = \rho,$$

$$P(Y = 1) = 1 - P(Y = -1) = \lambda,$$

and  $\rho$  and  $\lambda$  may be unknown, and so we regard them as random variables.

After observing  $X$  or  $Y$  he may choose to observe the same variable for the next bet or switch and observe the other variable. Which variable he chooses at the  $(n + 1)$ st stage may depend on the first  $n$  choices and observations. Using traditional bandit problem terminology, an observation of  $X$  is "a pull of the right arm  $R$ " and of  $Y$  is "a pull of the left arm  $L$ ."

DEFINITION 1.1. A strategy, sometimes denoted by  $\tau$ , is a function that assigns the symbol  $R$  or  $L$  to each finite history of bets and resulting observations.

The variable assigned by a strategy denotes the bet to be made next when following that strategy. Let  $Z_n$  denote the result of the  $n$ th bet and let  $S_n = \sum_{i=1}^n Z_i$ . The gambler's fortune at the  $(n + 1)$ st stage is

$$(1.1) \quad k_n = k + S_n.$$

If  $G < \infty$ , then the objective is to reach  $G$  before  $g$ . Any strategy that maximizes the probability that there exists an  $n$  such that

$$k_1 > g, \dots, k_{n-1} > g, k_n = G$$

will be called optimal. If  $G = \infty$  then any strategy that maximizes the

probability that  $k_n > g$  for all  $n$  and  $k_n \rightarrow \infty$  will be called optimal.

Any strategy that comes within  $\epsilon$  of the maximum, or supremum in case the maximum does not exist, will be called  $\epsilon$ -optimal. Evidently, if  $k_n = g$  or  $k_n = G$  for some  $n$ , the same is true for all larger  $n$ .

Let  $X_i$  and  $Y_j$  denote, respectively, the observations on the  $i$ th bet on  $\mathcal{R}$  and the  $j$ th bet on  $\mathcal{L}$ . For convenience we assume that the sequences  $(X_1, X_2, \dots)$  and  $(Y_1, Y_2, \dots)$  are non-terminating whether or not infinitely many bets are made on each variable.

The two sequences  $(X_1, X_2, \dots)$  and  $(Y_1, Y_2, \dots)$  are assumed to be independent of each other. Also, given the random variables  $\rho$  and  $\lambda$ ,  $(X_1, X_2, \dots)$  and  $(Y_1, Y_2, \dots)$  are sequences of independent random variables. Therefore, the unconditional finite-dimensional distributions of  $(X_1, X_2, \dots)$  and  $(Y_1, Y_2, \dots)$  are invariant under permutations of the subscripts; that is, each is a sequence of exchangeable random variables. In the present part, Part I, of this paper,  $\lambda$  will be assumed known with  $0 < \lambda < 1$ . Under this assumption,  $(Y_1, Y_2, \dots)$  is a sequence of independent random variables. The distribution  $R$  of  $\rho$  will be arbitrary. In Part II of this paper <sup>(Berry and Fristedt 1980)</sup>,  $\Lambda(\rho, \lambda)$  is known to be either  $(\alpha, \beta)$  or  $(\beta, \alpha)$ , where  $\alpha$  and  $\beta$  are arbitrary but known.

When possible, the notation used in this paper is consistent with that of Berry and Fristedt (1979), which deals with classical bandit problems. While many of the aspects of the current problem are different from those of classical bandit problems, some are the same. In particular, in both problems the desire to pull  $\mathcal{R}$  or  $\mathcal{L}$  depends on the probability of immediate success and on the possibility of gaining information about the variables  $X$  and  $Y$  which will aid future decisions.

In both problems these criteria may be conflicting, that is, suggesting different strategies. While it is not possible to consider these two criteria separately in either problem, the interrelationship is more complicated in the current problem.

For fortune  $k \in (g, G)$ , distribution  $R$  on  $[0,1]$ , and  $\lambda \in (0,1)$ , let  $U(k; R, \lambda)$  denote the maximum over the set of strategies, or supremum in case the maximum does not exist, of the probability of approaching fortune  $G$  (thereby reaching  $G$  when  $G$  is finite) before approaching  $g$ . The maximal probability of approaching  $G$  is at least that of the strategy which observes the same variable at each stage, and is no larger than if  $\rho$  is also known at the outset. That is,

$$(1.2) \quad u(k, \lambda) \vee Eu(k, \rho) \leq U(k; R, \lambda) \leq E[u(k, \rho) \vee u(k, \lambda)],$$

where  $u$  is the probability of success in the gambler's ruin problem. Expressions for this latter probability are well-known and are given below.

<u>Goals</u>	<u><math>u(k, p)</math></u>	<u>Restrictions</u>
$-\infty < g < G < \infty$	$\frac{1 - \left(\frac{1-p}{p}\right)^{k-g}}{1 - \left(\frac{1-p}{p}\right)^{G-g}}$	$p \neq 1/2$
$-\infty < g < G = \infty$	$\frac{k-g}{G-g}$	$p = 1/2$
$-\infty = g < G < \infty$	$1 - \left(\frac{1-p}{p}\right)^{k-g}$	$p \leq 1/2$ $p > 1/2$
$-\infty = g < G < \infty$	$\left(\frac{p}{1-p}\right)^{G-k}$	$p < 1/2$
	1	$p \geq 1/2$

We cannot find optimal strategies for the problem as stated in its full generality. We confine our attention mainly to special cases of  $(R, \lambda)$ . As we consider different cases, or even different instances within a single case, the reader will notice major qualitative differences in the optimal strategies. Some optimal strategies are counterintuitive. For example, it may be uniquely optimal to bet on  $Y$  (i.e., to pull  $\mathcal{E}$ ) even though it is practically certain that  $\rho > \lambda$ . The next example is especially easy and gives one instance of this phenomenon.

EXAMPLE 1.1. Suppose  $\lambda$  is .5 and

$$1 - R(\{0\}) = R(\{1\}) = r > 0.$$

Suppose both goals  $g$  and  $G$  are finite so we can take  $g = 0$ . We will abbreviate  $U(k; R, \lambda)$ , the maximal probability of reaching  $G$ , to  $U(k, r)$ . Again, there is no information to be gained by pulling  $\mathcal{E}$  since  $\lambda$  is known; that is,  $r$  is not affected by pulling  $\mathcal{E}$ . On the other hand, complete information is gained with probability one by a single pull of  $\mathcal{R}$ . It is reasonable, therefore, to expect that if  $r$  is large enough it is optimal to pull  $\mathcal{R}$  first and then proceed with complete knowledge. Actually, the optimal strategy is quite the opposite. If  $k = 2, \dots, G-1$ , it is optimal (uniquely unless  $r = 1$ ) to pull  $\mathcal{E}$  for all  $r$ , and it is optimal to pull  $\mathcal{R}$  if  $k = 1$ , for  $r \geq 1/G$  and  $\mathcal{E}$  for  $r < 1/G$ . Of course, if  $\mathcal{R}$  is ever pulled, then it is optimal (though not necessarily uniquely) to pull  $\mathcal{R}$  forever (until reaching  $G$ ) if it succeeds. Let  $\tau_0$  denote the strategy claimed to be optimal in these statements.

Let  $U_0(k, r)$  denote the probability of reaching  $G$  using  $\tau_0$  when the current fortune is  $k$ . Then

$$U_0(k, r) = \frac{k-1}{G-1} + \frac{G-k}{G-1} r, \quad \text{if } r \geq 1/G$$

$$= k/G, \quad \text{if } r < 1/G.$$

We show that  $U = U_0$  by showing that  $U_0$  is excessive (Corollary 2.1 of Section 2, or Dubins and Savage 1976, Theorem 2.12.1); that is, the expected value of  $U_0$  under either initial choice is not greater than  $U_0$ . It follows (Theorem 2.1 of Section 2) that the initial choice in question is optimal when and only when this expectation actually equals  $U_0$ . We prove that  $U_0$  is excessive by considering two cases that are determined by the arm pulled initially when following  $\tau_0$ .

According to  $\tau_0$ ,  $\mathcal{L}$  should be pulled when  $r < 1/G$  or both  $k > 1$  and  $r < 1$ . For such a fortune  $k$  and probability  $r$  that  $\rho = 1$ , the expected value of  $U_0$  under a pull of  $\mathcal{L}$  is simply  $U_0(k, r)$  from the definition of  $\tau_0$ . The expected value of  $U_0$  under a single pull of  $\mathcal{R}$  equals

$$rU_0(k+1, 1) + (1-r)U_0(k-1, 0)$$

$$= r + (1-r)(k-1)/G.$$

Subtracting this from  $U_0(k, r)$  gives

$$\left\{ \begin{array}{ll} (1-r) \frac{k-1}{G(G-1)} & \text{if } r \geq 1/G, \\ \frac{1}{G} - r + \frac{r}{G}(k-1) & \text{if } r < 1/G; \end{array} \right.$$

both of which are  $\geq 0$  (in fact  $> 0$  in the case under consideration) as desired.

In the remaining cases,  $\mathcal{R}$  is pulled when following  $\tau_0$ ; accordingly, we suppose that either  $k = 1$  and  $r \geq 1/G$  or  $k > 1$  and  $r = 1$ .

For such a fortune and probability that  $\rho = 1$  the expected value of  $U_0$  under a pull of  $\mathcal{R}$  is  $U_0(k,r)$  from the definition of  $\tau_0$ . The expected value of  $U_0$  under a single pull of  $\mathcal{L}$  equals

$$\frac{1}{2} U_0(2,r) = \frac{1}{2(G-1)} + \frac{G-2}{2(G-1)} r$$

if  $k = 1$ , and it equals 1 if  $k > 1$  and  $r = 1$ . Subtracting this from  $U_0(k,r)$  gives  $(rG-1)/(2G-2)$  if  $k = 1$  and 0 if  $k > 1$  and  $r = 1$ ; both of which are  $\geq 0$  for the cases at hand (in fact  $> 0$  unless  $k = 1$  and  $r = 1/G$  or  $k > 1$  and  $r = 1$ ) as desired.

The claim that  $\tau_0$  is optimal is established. Further, the only deviations from  $\tau_0$  that are optimal occur when equality holds in the excessivity calculations. So, pulling  $\mathcal{L}$  and pulling  $\mathcal{R}$  are both optimal when  $k = 1$  and  $r = 1/G$  or when  $k > 1$  and  $r = 1$ ; otherwise, the arm indicated by  $\tau_0$  is uniquely optimal.  $\square$

Two-armed bandits with a goal are motivated by the following medical treatment decision problem. Two treatments are available for use on a patient with a known level of a disease; one of the treatments may be "no treatment". Between applications the patient's condition either improves a little or worsens a little, independently of previous treatments and responses. Assume that levels can be quantified so that each new level is the old level plus or minus one. Furthermore, there is a level  $g$ , which is below the patient's current level, at which he can no longer be cured, and a higher level  $G$ , at which the patient no longer has the disease. Then the problem we have posed may be applicable. Perhaps the least realistic assumption we make is that it is only eventual cure or not that matters, and not the length of time that the patient is in a diseased state. Still, there are diseases

for which standard medical practice is similar to an optimal strategy in our problem for certain kinds of information. In particular, there are diseases in which treatments that are risky, having unknown effectiveness, are used only when the disease is in very advanced stages. As we saw in Example 1.1, some optimal strategies have this property. In that example, the right arm  $\alpha$  is a "risky treatment." On the other hand, there are diseases in which a risky treatment is only used when the patient is relatively healthy. We shall encounter situations in Section 5 in which optimal strategies have this property.

As far as we know, the problem considered here has not been studied previously, but it is a variation of a number of problems. Articles dealing with classical bandit problems--in which the expected number of successes (possibly discounted) is to be maximized--include those by Berry (1972) and Berry and Fristedt (1979). Related work in the theory of gambling includes the fundamental book of Dubins and Savage (1976), and (Berry, Heath, and Sudderth 1974). In a problem used as motivation by Dubins and Savage (1976), bets with known win probability are made with the objective of reaching a goal before going broke. The problem is to decide what sequence of stakes to use, depending on the sequence of fortunes. Berry, Heath, and Sudderth (1974) consider a generalization in which the win probability is itself a random variable. In the current problem there is a second sequence of bets available, but now only unit bets are allowed.

In Section 2 we introduce some necessary notation and terminology. In Section 3 we develop some general theory, considering the question of existence and, to a certain extent, the nature of optimal strategies.

In Section 4 we find optimal strategies when  $\rho$  is known to be 1 if it is greater than  $\lambda$ . In Section 5 we find optimal strategies when  $\rho$  is known to be 0 if it is less than  $\lambda$ . Finally, in Section 6 we consider the case in which  $G = \infty$  and  $R$  is a two-point distribution, the two points being symmetric about .5.

2. Preliminaries. The discussion in the previous section focusses on distinct initial fortunes  $k$  and distributions  $R$  of  $\rho$  without regard to their interconnections. In this section we knit these individual problems into one. Our purposes will be well-served by the following definition.

DEFINITION 2.1. For fixed  $\lambda$ ,  $g$ , and  $G$ , a scheme  $\Psi$  is a function  $\Psi(k,R)$  which assigns either of the symbols  $\mathcal{R}$  or  $\mathcal{L}$  to each  $k \in (g,G)$  and distribution  $R$  of  $\rho$ .

To connect the concepts of scheme and strategy, as defined in Definition 1.1, we need some further notation.

Just as we defined  $k_n$  in (1.1) to be the fortune at stage  $n+1$ , we let  $R_n$  be the conditional distribution of  $\rho$  at stage  $n+1$ . Specifically, if  $s$  successes and  $f$  failures have been obtained on  $\mathcal{R}$  in the first  $n$  stages, then  $R_n = \sigma^s \varphi^f R$ , where each  $\sigma^s \varphi^f R$  is absolutely continuous with respect to  $R$  and

$$\frac{d\sigma^s \varphi^f R}{dR}(x) = \frac{x^s (1-x)^f}{\int u^s (1-u)^f R(du)}.$$

(In case the support of  $R$  is contained in  $\{0,1\}$ , written  $\text{supp } R \subset \{0,1\}$ , some values of  $(s,f)$  are not possible.) Stage by stage considerations yield the characterization:  $R_0 = R$  and

$$\begin{aligned} R_n &= \sigma R_{n-1} && \text{if } \mathcal{R} \text{ is pulled at stage } n \text{ and } Z_n = 1 \\ &= \varphi R_{n-1} && \text{if } \mathcal{R} \text{ is pulled at stage } n \text{ and } Z_n = -1 \\ &= R_{n-1} && \text{if } \mathcal{L} \text{ is pulled at stage } n. \end{aligned}$$

DEFINITION 2.2. The strategies induced by the scheme  $\Psi$  are those that require a pull of  $\Psi(k_n, R_n)$  at the  $(n+1)$ st stage, for  $n = 0, 1, \dots$ .

In Example 1.1 we restricted attention to strategies induced by a scheme, that is, strategies that specify the arm to be pulled at each stage depending only on the current fortune and current distribution of  $\rho$ , and not otherwise on the history of observations. While there may be optimal strategies outside this class, we shall see in Theorems 3.1 and 3.2 that there are always strategies inside the class that perform at least as well as those outside.

For the consideration of schemes we shall require a notation that reflects the simultaneous consideration of a variety of initial fortunes  $k$  and distributions  $R$ . For a strategy  $\tau$ , let  $E_{\tau}^{(k,R)}$  and  $P_{\tau}^{(k,R)}$  (or  $P_{\psi}^{(k,R)}$  when  $\tau$  is induced by the scheme  $\psi$ ) denote the expectation and probability starting at  $(k,R)$ . For instance,

$$U(k;R,\lambda) = \sup_{\tau} P_{\tau}^{(k,R)}(k_n \rightarrow G).$$

For situations in which expectation depends on the strategy  $\tau$  only through the arm pulled initially, we shall use the name of that arm as the subscript. In this notation, we can express what is called the "functional equation" in sequential decision theory:

$$(2.1) \quad U(k;R,\lambda) = E_R^{(k,R)} U(k_1;R_1,\lambda) \vee E_g^{(k,R)} U(k_1;R_1,\lambda).$$

When convenient, we shall drop any of the superscripts or subscripts from the  $E$  and  $P$  notations. For instance, the probability of a second success with  $g$  given one success can be written in a variety of ways:

$$\begin{aligned} E^{gR}(\rho) &= \int x (\sigma R)(dx) = \frac{\int x^2 R(dx)}{\int x R(dx)} \\ &= \frac{E^R(\rho^2)}{E^R(\rho)} = \frac{E\rho^2}{E\rho}. \end{aligned}$$

An arm  $G$ , where  $G$  equals  $R$  or  $F$ , is called conserving at  $(k,R)$  if

$$E_G^{(k,R)} U(k_1; R_1, \lambda) = U(k; R, \lambda);$$

that is, if the maximum in (2.1) is taken on at arm  $G$ . Clearly, it is never optimal to pull an arm that is not conserving. However, as we shall see, it may not be optimal to always pull a conserving arm.

To carry over the concepts of conserving and optimal to schemes, we first need the following definition.

DEFINITION 2.3. Let  $\underline{R}$  be a family of distributions of  $\rho$ . We say that  $\underline{R}$  is closed under sampling whenever  $R \in \underline{R}$  implies  $\sigma^s \varphi^f R \in \underline{R}$  for all  $s$  and  $f$  that are possible under  $R$ . If  $R((0,1)) > 0$ , then all pairs of nonnegative integers  $(s,f)$  are possible.

DEFINITION 2.4. A scheme  $\Psi$  is called conserving for  $\underline{R}$ , a class of distributions that is closed under sampling, if  $\Psi(k,R)$  is conserving at  $(k,R)$  for each  $k \in (g,G)$  and  $R \in \underline{R}$ .

DEFINITION 2.5. Let  $\epsilon \geq 0$ . A scheme  $\Psi$  for which the induced strategy for each initial  $k \in (g,G)$  and each initial  $R \in \underline{R}$ , a class of distributions closed under sampling, is  $\epsilon$ -optimal is called an  $\epsilon$ -optimal scheme for  $\underline{R}$ . A 0-optimal scheme for  $\underline{R}$  is called an optimal scheme for  $\underline{R}$ .

An arm  $G$  is called optimal at  $(k,R)$  if  $G = \Psi(k,R)$  for some scheme  $\Psi$  that is optimal for some class containing  $R$ .

While we give an elementary proof of the following theorem, it can be viewed as a special case of Theorem 14 of Sudderth (1972), which in turn is a generalization of Theorem 3.5.1 of Dubins and Savage (1976). In case  $G < \infty$ , our theorem is also a special case of the Dubins and Savage result.

THEOREM 2.1. A conserving scheme  $\Psi$  for a class  $\underline{R}$ , closed under sampling, is optimal for  $\underline{R}$  if, for each  $k \in (g, G)$  and  $R \in \underline{R}$ ,

$$(2.2) \quad P_{\Psi}^{(k, R)}(U(k_n; R_n, \lambda) \rightarrow 0 | k_n \rightarrow G) = 1 .$$

PROOF. Suppose  $\Psi$  is conserving for  $\underline{R}$ . By induction,

$$U(k; R, \lambda) = E_{\Psi}^{(k, R)} U(k_n; R_n, \lambda) .$$

Let  $n \rightarrow \infty$  to obtain

$$U(k; R, \lambda) \leq P_{\Psi}^{(k, R)}(k_n \rightarrow G) + P_{\Psi}^{(k, R)}(k_n \rightarrow G) \text{ess sup}_{[k_n \rightarrow G]} \limsup_{n \rightarrow \infty} U(k_n; R_n, \lambda) ,$$

where  $\text{ess sup}_{[k_n \rightarrow G]}$  denotes the essential supremum for the measure

$P_{\Psi}^{(k, R)}$  over the event  $[k_n \rightarrow G]$ . ■

We now make precise the concept of excessivity used in Example 1.1.

DEFINITION 2.6. Let  $W(k, R)$  be a real-valued function defined for  $g-1 < k < G+1$  and  $R \in \underline{R}$ , a class closed under sampling. The function  $W$  is excessive for  $\underline{R}$  if, for each arm  $G$  and each  $(k, R) \in (g-1, G+1) \times \underline{R}$ ,

$$W(k, R) \geq E_G^{(k, R)} W(k_1, R_1) .$$

The following theorem enables one to check that a particular scheme is optimal. While we give an elementary proof of the theorem, it can be viewed as a special case of Theorem 1 of Dubins and Sudderth (1977), which in turn is a generalization of Theorem 2.12.1 of Dubins and Savage (1976).

THEOREM 2.2. Suppose that, for some scheme  $\Psi$ , the function

$$(k, R) \mapsto P_{\Psi}^{(k, R)}(k_n \rightarrow G)$$

is excessive for some class  $\underline{R}$  that is closed under sampling. If, in addition,

$$(2.3) \quad P_{\Psi}^{(k,R)}(k_n \rightarrow G) \rightarrow 1 \text{ as } k \rightarrow G, \text{ uniformly for } R \in \underline{R},$$

then  $\Psi$  is optimal for  $\underline{R}$ .

PROOF. Fix  $(k,R) \in (g,G) \times \underline{R}$  and let  $\tau$  be an arbitrary strategy for starting at  $(k,R)$ . Let  $\tau_m$  denote the strategy: use  $\tau$  at stages 1, ..., m and thereafter use the strategy induced by  $\Psi$ . By induction and the assumed excessivity,

$$(2.4) \quad P_{\tau_m}^{(k,R)}(k_n \rightarrow G) \leq P_{\Psi}^{(k,R)}(k_n \rightarrow G).$$

On the other hand,

$$(2.5) \quad \begin{aligned} & \liminf_{m \rightarrow \infty} P_{\tau_m}^{(k,R)}(k_n \rightarrow G) \\ &= \liminf_{m \rightarrow \infty} P_{\tau}^{(k,R)}(P_{\Psi}^{(k_m, R_m)}(k_n \rightarrow G) | (k_m, R_m)) \\ &> P_{\tau}^{(k,R)}(k_m \rightarrow G) \inf_{[k_m \rightarrow G]} \liminf_{m \rightarrow \infty} P_{\Psi}^{(k_m, R_m)}(k_n \rightarrow \infty). \end{aligned}$$

By (2.3) the right side of (2.5) equals  $P_{\tau}(k_m \rightarrow G)$  which, by (2.4) and (2.5) is less than or equal to  $P_{\Psi}(k_n \rightarrow G)$ . ■

For  $G < \infty$ , Theorem 2.2 is a special case of Theorem 2.12.1 of Dubins and Savage (1976), and since (2.3) is trivially satisfied, simplifies as follows.

COROLLARY 2.1. Suppose  $G < \infty$ . If  $(k,R) \rightarrow P_{\Psi}^{(k,R)}(k_n \rightarrow G)$  is excessive for  $\underline{R}$ , then  $\Psi$  is optimal for  $\underline{R}$ .

3. Optimal Schemes; General Theory. In this section we prove two theorems dealing with the existence of optimal schemes and one giving a qualitative property of optimal schemes.

THEOREM 3.1. There is an optimal scheme for the class of all distributions of  $\rho$  if  $g > -\infty$  or  $G = \infty$  or  $\lambda \geq .5$ .

PROOF. (i) Suppose both  $G = \infty$  and  $\lambda < .5$ . Then the only chance of approaching  $\infty$  is if  $\rho > .5$ ; therefore, it is optimal to pull  $R$  at each stage.

(ii) Suppose  $g = -\infty$ ,  $\lambda \geq .5$ , and either  $G < \infty$  or  $\lambda \neq .5$ . Then it is clearly optimal to pull  $R$  at each stage.

(iii) Finally, suppose  $g > -\infty$  and either  $G < \infty$  or  $\lambda > .5$ . Since we are assuming  $\lambda < 1$  throughout (thereby ruling out the possibility of oscillating forever), the event  $[k_n \neq G]$  differs by a null event from the event  $[k_n \rightarrow g]$  (that is, the event  $[k_n = g$  for some  $n]$ ). Thus, since  $g > -\infty$ , condition (2.2) is satisfied for every scheme  $\psi$ . It follows from Theorem 2.1 that every conserving scheme is optimal. Since there exists at least one conserving scheme, there is at least one optimal scheme. ■

The only case not covered in Theorem 3.1 will be considered in Theorem 3.2. For one circumstance arising in that theorem we shall need the following lemma.

LEMMA 3.1. Suppose that  $\lambda < \lambda + \delta < .5 - \delta$  and  $g = -\infty$ . Then there exists  $\gamma > 0$  such that

$$\frac{R([0, \lambda + \delta])}{R([\lambda + \delta, .5 - \delta])} < \gamma$$

implies  $E_{\mathcal{L}}^{(k, R)}(U(k_1; R_1, \lambda)) < E_{\mathcal{R}}^{(k, R)}(U(k_1; R_1, \lambda))$  for every  $k$ .

PROOF. From (1.2), we obtain

$$E_{\mathcal{F}}^{(k,R)} U(k_1; R_1, \lambda) \leq \lambda E^R [u(k+1, \rho) \vee u(k+1, \lambda)] + (1-\lambda) E^R [u(k-1, \rho) \vee u(k-1, \lambda)]$$

and

$$\begin{aligned} E_{\mathcal{R}}^{(k,R)} U(k_1; R_1, \lambda) &\geq E^R u(k, \rho) \\ &= E[\rho u(k+1, \rho) + (1-\rho) u(k-1, \rho)]. \end{aligned}$$

Therefore

$$\begin{aligned} &E_{\mathcal{R}}^{(k,R)} U(k_1; R_1, \lambda) - E_{\mathcal{F}}^{(k,R)} U(k_1; R_1, \lambda) \\ &\geq E^R \left[ \mathbb{1}_{[0, \lambda+\delta]}^{(\rho)+1} \mathbb{1}_{[\lambda+\delta, .5-\delta]}^{(\rho)+1} \mathbb{1}_{(.5-\delta, 1]}^{(\rho)} \{ \rho u(k+1, \rho) \right. \\ &\quad \left. + (1-\rho) u(k-1, \rho) - \lambda [u(k+1, \rho) \vee u(k+1, \lambda)] - (1-\lambda) [u(k-1, \rho) \vee u(k-1, \lambda)] \right] \\ &= I_1 + I_2 + I_3, \end{aligned}$$

where  $I_1$ ,  $I_2$ , and  $I_3$  denote the terms involving the respective indicator functions.

Clearly,

$$I_1 \geq E^R \left[ \mathbb{1}_{[0, \lambda+\delta]}^{(\rho)} \{ -\lambda u(k+1, \lambda+\delta) - (1-\lambda) u(k-1, \lambda+\delta) \} \right] \geq -R([0, \lambda+\delta]) u(k+1, \lambda+\delta).$$

We also have

$$\begin{aligned} I_2 &= E^R \left[ \mathbb{1}_{[\lambda+\delta, .5-\delta]}^{(\rho)} (\rho - \lambda) \{ u(k+1, \rho) - u(k-1, \rho) \} \right] \\ &\geq \delta R([\lambda+\delta, .5-\delta]) \inf \{ u(k+1, p) - u(k-1, p) : \lambda + \delta \leq p \leq .5 - \delta \} \end{aligned}$$

and, by a similar calculation,  $I_3 \geq 0$ . Since, for  $p \in [\lambda+\delta, .5-\delta]$ ,

$$\begin{aligned} u(k+1, p) - u(k-1, p) &= u(k+1, p) - \left( \frac{p}{1-p} \right)^2 u(k+1, p) \\ &\geq u(k+1, p) \left[ 1 - \left( \frac{.5-\delta}{.5+\delta} \right)^2 \right] \geq u(k+1, \lambda+\delta) \left[ 1 - \left( \frac{.5-\delta}{.5+\delta} \right)^2 \right], \end{aligned}$$

$I_1 + I_2 + I_3 > 0$  for  $R([0, \lambda+\delta])/R([\lambda+\delta, .5-\delta]) < \gamma$ , where  $\gamma = \delta [1 - (.5-\delta)^2 / (.5+\delta)^2]$ . ■

THEOREM 3.2. Suppose  $g = -\infty$ ,  $G < \infty$ , and  $\lambda < .5$ . For each  $\epsilon > 0$  there is an  $\epsilon$ -optimal scheme for the class of all distributions of  $\rho$ . There is an optimal scheme for  $R_3 \cup R_4 \cup R_5$  where

$$R_3 = \{R: R((\lambda, .5)) > 0\},$$

$$R_4 = \{R: R([0, \lambda]) = 1\}.$$

$$R_5 = \{R: R([\lambda, 1]) = 1\}.$$

There is no optimal strategy for any initial  $(k, R)$  for  $R$  belonging to the complement of  $R_3 \cup R_4 \cup R_5$ .

REMARKS. Each of the classes  $R_3$ ,  $R_4$ , and  $R_5$  is closed under sampling, and therefore, so is their union. The last sentence in the theorem implies that  $R_3 \cup R_4 \cup R_5$  is the largest class closed under sampling for which there is an optimal scheme.

PROOF OF THEOREM 3.2. Let  $R_6$  denote the complement of  $R_3 \cup R_4 \cup R_5$ , so

$$R_6 = \{R: R((\lambda, .5)) = 0, R([0, \lambda]) > 0, R([\lambda, 1]) > 0\}.$$

We shall define a scheme  $\Psi$  and show that it induces an optimal strategy for any  $R \in R_3 \cup R_4 \cup R_5$  and an  $\epsilon$ -optimal strategy for any  $R \in R_6$ .

Define  $\Psi(k, R) = \mathcal{L}$  for  $R \in R_4$  and all  $k$ . Define  $\Psi(k, R) = \mathcal{R}$  for  $R \in R_5 - R_4$  and all  $k$ . Clearly, as so far defined,  $\Psi$  is an optimal scheme for  $R_4 \cup R_5$ .

For  $R \in R_6$  define

$$\Psi(k, R) = \begin{cases} \mathcal{L} \\ \mathcal{R} \end{cases} \text{ for } k \begin{cases} > \\ < \end{cases} k^* = \left[ G - \frac{\log \epsilon}{\log[\lambda/(1-\lambda)]} \right],$$

the largest integer less than or equal to  $G - \log \epsilon / \log[\lambda/(1-\lambda)]$ .

Since  $R_4 \cup R_5 \cup R_6$  is closed under sampling, we can calculate, for  $R \in R_6$ , that the probability of reaching  $G$  when starting at  $k$  and following  $\Psi$  is at least

$$\begin{aligned} & 1 - \frac{1 - [\lambda/(1-\lambda)]^{G-(k \vee k^*)}}{1 - [\lambda/(1-\lambda)]^{G-k^*}} R([0, \lambda]) \\ & \geq 1 - [1 - (\frac{\lambda}{1-\lambda})^{G-(k \vee k^*)} + (\frac{\lambda}{1-\lambda})^{G-k^*}] R([0, \lambda]) \\ & \geq 1 - [1 - (\frac{\lambda}{1-\lambda})^{G-k} + \epsilon] R([0, \lambda]) . \end{aligned}$$

Therefore,

$$(3.1) \quad U(k; R, \lambda) \geq 1 - R([0, \lambda]) (1 - [\lambda/(1-\lambda)]^{G-k}).$$

The lower bound in (3.1) is precisely the upper bound given in (1.2), so that for  $R \in R_6$ ,

$$U(k; R, \lambda) = 1 - R([0, \lambda]) (1 - [\lambda/(1-\lambda)]^{G-k}),$$

which means that  $\Psi$  is  $\epsilon$ -optimal. The supremum  $U$  can only be attained if  $R$  is pulled eventually (and often) when  $\rho \geq .5$  and if  $\rho$  is pulled exclusively when  $\rho < \lambda$ . Since, for  $R \in R_6$ , no strategy fulfills both these requirements with probability one, there is no optimal strategy for  $R \in R_6$ , the complement of  $R_3 \cup R_4 \cup R_5$ .

It remains to consider  $R_3$ , a class closed under sampling. Define  $\Psi$  to be conserving for  $R_3$ . For each  $R \in R_3$  there exists a  $k^*$  such that: for  $k < k^*$

$$\begin{aligned} & E_{\rho}^{(k, R)} U(k_1; R_1, \lambda) \\ & \leq \lambda E^R [u(k+1, \lambda) \vee u(k+1, \rho)] + (1-\lambda) E^R [u(k-1, \lambda) \vee u(k-1, \rho)] \\ & < E^R [\rho u(k+1, \rho) + (1-\rho) u(k-1, \rho)] \end{aligned}$$

$$\begin{aligned}
&= E^R u(k, \rho) \\
&\leq E_{\theta}^{(k, R)} U(k_1; R_1, \lambda)
\end{aligned}$$

and, hence,  $\Psi(k, R) = \theta$ . Therefore, for any  $(k, R) \in (-\infty, G) \times R_3$ ,

$$(3.2) \quad P_{\Psi}^{(k, R)}(\Psi(k_n, R_n) = \theta \text{ i.o. } | k_n \rightarrow -\infty) = 1.$$

By the strong law of large numbers

$$(3.3) \quad P_{\Psi}(s_n / (s_n + f_n) \rightarrow \rho | k_n \rightarrow -\infty) = 1,$$

where  $s_n$  and  $f_n$  denote the numbers of successes and failures on  $\theta$  for stages 1 through  $n$ .

Choose  $\delta > 0$ , depending on  $R \in R_3$ , so that  $R([0, \lambda + \delta], [\lambda + \delta, .5 - \delta]) > 0$ .

An easy calculation starting with (3.3) gives

$$P_{\Psi}\left(\frac{R_n([0, \lambda + \delta])}{R_n([\lambda + \delta, .5 - \delta])} \rightarrow 0 | k_n \rightarrow -\infty, \rho \geq .5\right) = 1.$$

Now, by Lemma 3.1,

$$P_{\Psi}(\Psi(k_n, R_n) = \xi \text{ only finitely often } | k_n \rightarrow -\infty, \rho \geq .5) = 1$$

which implies

$$P_{\Psi}(k_n \rightarrow G | k_n \rightarrow -\infty, \rho \geq .5) = 1$$

and, thus,

$$(3.4) \quad P_{\Psi}(\rho \geq .5 | k_n \rightarrow -\infty) = 0.$$

From (3.2) and the weak law of large numbers we obtain

$$P_{\Psi}(R_n \rightarrow \delta(\rho) | k_n \rightarrow -\infty) = 1,$$

where  $\delta(x)$  denotes the probability measure with mass 1 at  $x$ . Hence, using (1.2),

$$P_{\Psi}(U(k_n; R_n, \lambda) \rightarrow 0 | k_n \rightarrow -\infty, \rho < .5) = 1.$$

By (3.4) we can drop the condition  $\rho < .5$  and because we are assuming  $\lambda > 0$  throughout we can replace  $k_n \rightarrow -\infty$  by  $k_n \nearrow G$ . Hence, (2.2) holds, and so Theorem 2.1 implies that  $\Psi$  is optimal for  $R_3$ .  $\square$

Because there is no gain in information when  $\mathcal{F}$  is pulled, it seems unlikely for there to exist a fortune  $k$  at which  $\mathcal{F}$  is uniquely optimal and yet  $\mathcal{R}$  is optimal at fortunes  $k-1$  and  $k+1$ . The next theorem verifies this, saying that the set of fortunes for which  $\mathcal{R}$  is optimal is contiguous.

**THEOREM 3.3.** The optimal and  $\varepsilon$ -optimal schemes of Theorems 3.1 and 3.2 may be chosen from among the set of schemes  $\Psi$  for which  $\{k: \Psi(k, R) = \mathcal{R}\}$  is an interval. Furthermore, if  $G - g = \infty$ , the schemes of Theorems 3.1 and 3.2 may be chosen from among the set of schemes  $\Psi$  for which  $\{k: \Psi(k, R) = \mathcal{R}\}$  is empty or an unbounded interval.

**PROOF.** Let  $\Psi(k, R) = \mathcal{R}$  if  $\mathcal{R}$  is conserving at  $(k, R)$  and  $\mathcal{F}$  otherwise. By Definition 2.4,  $\Psi$  is conserving for the class of all distributions of  $\rho$ . Using a well-known argument, we shall show by contradiction that  $\{k: \Psi(k, R) = \mathcal{R}\}$  is an interval.

Suppose for a certain  $R$  and  $k' \leq k'' - 2$  that  $\Psi(k', R) = \mathcal{R}$ ,  $\Psi(k'', R) = \mathcal{R}$ , and  $\Psi(k, R) = \mathcal{F}$  for all  $k$  with  $k' < k < k''$ . For this certain  $R$  and some  $k \in (k', k'')$  let  $\tau$  denote the strategy: pull  $\mathcal{R}$  at stage 1, then pull  $\mathcal{F}$  until either fortune  $k' + 1$  or fortune  $k'' + 1$  is reached in case  $Z_1 = 1$  or until either fortune  $k' - 1$  or fortune  $k'' - 1$  is reached in case  $Z_1 = -1$ , and thereafter follow  $\Psi$ . Then

$$P_{\tau}^{(k, R)}(k_n \rightarrow G) = P_{\Psi}^{(k, R)}(k_n \rightarrow G).$$

Therefore,  $\varrho$  is conserving at  $(k,R)$  and, hence,  $\Psi(k,R) = \varrho$ , the desired contradiction. When  $G - g = \infty$ , similar reasoning shows that  $\{k: \Psi(k,R) = \varrho\}$  is either empty or an unbounded interval.

Therefore,  $\Psi$  is a conserving scheme with the desired characteristic.

For  $g > -\infty$  and either  $G < \infty$  or  $\lambda > .5$ , Theorem 2.1 together with the fact that the difference between  $[k_n \rightarrow G]$  and  $[k_n \rightarrow g]$  is a null set implies that  $\Psi$  is optimal for the class of all distributions of  $\rho$ . For  $g = -\infty$ ,  $G < \infty$ ,  $\lambda < .5$ , the proof of Theorem 3.2 shows every conserving scheme for  $R_3$ , in particular  $\Psi$ , to be optimal for  $R_3$ .

For the remaining cases the schemes explicitly constructed in the proofs of Theorems 3.1 and 3.2 satisfy Theorem 3.3. ■

Suppose that we were to replace the objective we have used throughout by the objective: maximize the probability that  $k_n \rightarrow g$ . None of the conclusions in this paper would change except in case  $g = -\infty$ ,  $G = \infty$ , and  $\lambda = .5$ . With the change in objective,  $\Psi(k,R) \equiv g$  would become optimal in this case.

The remainder of this paper considers particular subclasses of the class of all distributions of  $\rho$ . The reader familiar with the emphasis on "stay with the winner" strategies in the literature will notice the various situations in which it is not optimal to pull an arm immediately after having obtained a success with that arm when following an optimal strategy.

4. Optimal Schemes for the Class  $R_{\lambda}$ . In this section we find all optimal schemes when  $\rho$  is known to be either 1 or no greater than  $\lambda$ , thereby generalizing Example 1.1. Define

$$R_{\lambda} = \{R: \text{supp } R \subset [0, \lambda] \cup \{1\}\},$$

a class of distributions that is closed under sampling. If  $R$  is in this class, then  $\text{supp } \varphi R \subset [0, \lambda]$ ; so the effectiveness of  $f$  is known to be at least that of  $R$  if ever  $R$  yields a failure. Of course, no finite number of successes on  $R$  is as revealing in this regard as a single failure.

For reasons discussed in the previous section, the case  $g = -\infty$  is rather trivial, so we exclude it from consideration. Also, if  $G = \infty$  and  $\lambda \leq .5$  then the problem is trivial; the only chance of winning is that  $\rho = 1$ , so  $U$  is just  $R(\{1\})$ . We therefore exclude this case as well.

Let  $v(\rho)$  denote the probability, given  $\rho$ , of approaching  $G$  rather than reaching  $g$  when starting at  $g + 1$  and using the strategy: pull  $R$  at  $g + 1$  and  $f$  at all other fortunes. Standard calculations yield:

$$\begin{aligned} v(\rho) &= \frac{\rho[1 - (1-\lambda)/\lambda]}{1-\rho + \rho[1-(1-\lambda)/\lambda]} \quad \text{if } \lambda > .5 \text{ and } G = \infty \\ &= \frac{\rho[1 - (1-\lambda)/\lambda]}{(1-\rho)[1 - ((1-\lambda)/\lambda)^{G-g-1}] + \rho[1-(1-\lambda)/\lambda]} \quad \text{if } \lambda \neq .5 \text{ and } G < \infty \\ &= \frac{\rho}{(1-\rho)(G-g-1) + \rho} \quad \text{if } \lambda = .5 \text{ and } G < \infty. \end{aligned}$$

The next theorem generalizes Example 1.1, which concerns the case  $\text{supp } R \subset \{0,1\}$  and  $G < \infty$ .

**THEOREM 4.1.** Assume that  $g > -\infty$  and either  $G < \infty$  or  $\lambda > .5$ .

Then

$$\Psi_1(k,R) = \begin{cases} R & \text{if } k = g+1 \text{ and } E^R v(\rho) > v(\lambda) \\ f & \text{otherwise} \end{cases}$$

is an optimal scheme for  $R_1$ .

**REMARK.** From the definition of  $u$  in Section 1, it is evident that

$$v(\lambda) = u(g+1, \lambda).$$

**PROOF OF THEOREM 4.1.** By calculations, which we omit, that are similar to those used in Example 1.1,

$$V_1(k,R) = P_{\Psi_1}^{(k,R)}(k_n \rightarrow G)$$

can be shown to be excessive. For  $G < \infty$  the optimality of  $\Psi_1$  follows from Corollary 2.1. For  $G = \infty$ , Theorem 2.2 applies; for, (2.3) is satisfied by  $\Psi_1$  since  $\lambda > .5$  and  $\Psi_1(k,R) = f$  for  $k > g+1$ .

One of the calculations mentioned in the above proof shows that

$$V_1(k,R) \geq E(\rho)V_1(k+1,OR) + E(1-\rho)V_1(k-1,OR)$$

for  $k > g+1$ , with equality if and only if  $G = \infty$  and  $E v(\rho) \geq v(\lambda)$ .

Keeping track of this and other situations in which equality holds when checking the excessivity of  $V_1$  allows us to list every conserving scheme. By Theorem 2.1 these are exactly the optimal schemes. Accordingly, we have the following theorem.

**THEOREM 4.2.** Assume that  $g > -\infty$  and either  $G < \infty$  or  $\lambda > .5$ .

A scheme  $\Psi$  is optimal for  $R_1$  if and only if

$$\begin{aligned}
V(k,R) &= R \text{ if } k = g+1 \text{ and } E^R v(\rho) > v(\lambda) \\
&= I \text{ if } k > g+1, E((\lambda,1)) < 1, \text{ and } G < \infty \\
&= I \text{ if } E^R v(\rho) < v(\lambda).
\end{aligned}$$

REMARKS. According to this theorem, there are many cases in which either arm can be pulled without losing optimality. For example, if  $k > g+1$ ,  $E v(\rho) \geq v(\lambda)$ , and  $G = \infty$ , then either arm can be pulled. This makes it clear that there are many optimal strategies that are not induced by any optimal schema. One such strategy in the case  $G = \infty$  makes the following minor modification of some strategy induced by  $\Psi_1$ : arm @ is used when (and if) fortune  $g+2$  is reached for the second time provided that  $E v(\rho) \geq v(\lambda)$ .

In the next section we consider a class of distributions which is at the other extreme from the class considered in this section.

5. Optimal Schemes for the Class  $R_2$ . In this section we investigate optimal schemes when  $\rho$  is known to be either 0 or at least as large as  $\lambda$ . Define

$$R_2 = \{R: \text{supp } R \subset \{0\} \cup [\lambda, 1], R((\lambda, 1)) > 0\},$$

a class of distributions that is closed under sampling. We have introduced the condition  $R((\lambda, 1)) > 0$  since R's for which  $R(\{0, \lambda, 1\}) = 1$  are contained in  $R_1$  and so are covered in the previous section.

For a distribution  $R \in R_2$ ,  $\text{supp } R \subset [\lambda, 1]$ , so the effectiveness of  $R$  is known to be at least that of  $g$  if ever  $R$  yields a success. In a sense  $R_2$  is at the opposite extreme from  $R_1$  since it is now a single success rather than a failure that settles the issue of which is the better arm. While the classes  $R_1$  and  $R_2$  are comparable in this regard, the matter of finding optimal schemes for  $R_2$  is much more complicated than for  $R_1$ .

We shall construct a scheme  $\Psi_2$  only on  $(g, G) \times R_2$  and prove it to be optimal for  $R_2$ , checking the ways a scheme may differ from  $\Psi_2$  and still be optimal for  $R_2$ . For "typical" R's belonging to  $R_2$  it will develop that  $\Psi(k, R)$  must equal  $\Psi_2(k, R)$  for all  $k$  in order that  $\Psi$  be optimal. Throughout we exclude the case  $g = -\infty$ ; our methods do not work in that case.

We first consider the relatively simple case  $G = \infty$  and then examine the modifications necessary when  $G < \infty$ .

Suppose  $G = \infty$  and, to exclude a trivial case discussed in Section 3, assume  $\lambda > .5$ . Propositions 5.1 and 5.2 allow us to construct the scheme  $\Psi_2$  recursively.

Proposition 5.3 makes it clear that, for any fixed  $(k, R)$ ,  $V_2(k, R)$  can be calculated in a finite number of steps. Then Proposition 5.2 and Theorem 5.1 imply that, for any fixed initial  $(k, R)$ , an optimal strategy (for all stages) and

$$P_{V_2}^{(k, R)}(k_n \rightarrow G)$$

both can be calculated in a finite number of steps. All optimal strategies for  $R_2$  are characterized in Theorem 5.1.

To define  $V_2$  we need the following notation:

$$(5.1) \quad S_0 = \{R \in R_2 : E^R(\rho | \rho \neq \lambda) \leq (2\lambda - 1)/\lambda\};$$

$$(5.2) \quad S_m = \{R \in R_2 : \forall R \in S_{m-1}, R \notin S_{m-1}\}, m = 1, 2, \dots;$$

$$(5.3) \quad S_\infty = \{R \in R_2 : R(\{0\}) = 0\}.$$

If  $R \in S_\infty$ , then

$$E^{\rho^m R}(\rho | \rho \neq \lambda) > \lambda > (2\lambda - 1)/\lambda$$

for all  $m$ , and if  $R \in R_2 - S_\infty$ , then

$$E^{\rho^m R}(\rho | \rho \neq \lambda) \rightarrow 0 \text{ as } m \rightarrow \infty.$$

Therefore, we have the next result.

PROPOSITION 5.1.  $\{S_0, S_1, \dots, S_\infty\}$  is a partition of  $R_2$ .

For easy reference we formally state the following obvious facts.

PROPOSITION 5.2. The classes  $S_\infty$  and  $S_\infty \cup \bigcup_{i=0}^m S_i$ ,  $m = 0, 1, 2, \dots$ , are closed under sampling. Moreover, for  $R \in S_m$ ,  $1 \leq m < \infty$ ,  $\forall R \in S_m$  and  $\forall R \in S_{m-1}$ .

Let

$$(5.4) \quad \begin{aligned} \Psi_2(k, R) &= g \quad \text{if } (k, R) \in (g, G) \times S_{-\infty} \\ &= f \quad \text{if } (k, R) \in (g, G) \times S_{-0}. \end{aligned}$$

Suppose  $\Psi_2(k, R)$  has been defined for  $R \in S_{-\infty} \cup \bigcup_{i=0}^{\infty} S_{-i}$  and all  $k$ .

Consider an  $R$  belonging to  $S_{-g+1}$ . For starting at  $(g+1, R)$  let

$\tau_j$  denote the strategy: pull  $f$  until fortune  $j$  is reached (if ever), then pull  $g$  once, and thereafter follow  $\Psi_2$ . For  $j < \infty$ , let

$$(5.5) \quad q_j(R) = P_{\tau_j}^{(g+1, R)}(k_n \rightarrow G)$$

which can be calculated using the portion of  $\Psi_2$  already constructed.

Let

$$(5.6) \quad a(R) = \inf\{k: q_k(R) = \sup_j q_j(R)\},$$

$$(5.7) \quad \begin{aligned} \Psi_2(k, R) &= g \quad \text{if } k < a(R) \\ &= g \quad \text{if } k \geq a(R). \end{aligned}$$

From Propositions 5.1 and 5.2 we see that we have completed a recursive definition of a scheme  $\Psi_2$ . The next result makes it clear that we can use (5.6) to calculate  $a(R)$  in a finite number of steps.

**PROPOSITION 5.3.** Suppose  $G = \infty$ . For  $R \in \bigcup_{m=1}^{\infty} S_{-m}$ , the sequence  $j \mapsto q_j(R)$  defined by (5.5) satisfies:

- (i)  $q_j(R) \rightarrow (2\lambda-1)/\lambda$  as  $j \rightarrow \infty$ ;
- (ii)  $q_j(R) > (2\lambda-1)/\lambda$  for some  $j$ ;
- (iii)  $q_j(R) < \frac{(2\lambda-1)/\lambda}{1 - (\frac{1-\lambda}{\lambda})^{j-g}}$  for all  $j$ .

**REMARKS.** By calculating, for  $R \in \bigcup_{m=1}^{\infty} S_{-m}$ , the numbers  $q_1(R), q_2(R), \dots$  in order, one will eventually find some  $j$  and

some  $y > (2\lambda-1)/\lambda$  for which  $y = q_j(R)$ . For  $j$  greater than or equal to some  $J$  it will be true that

$$\frac{(2\lambda-1)/\lambda}{1 - \left(\frac{1-\lambda}{\lambda}\right)^{j-g}} < y.$$

By Proposition 5.3, only  $q_j(R)$  for  $j < J$  need be examined in order to calculate  $a(R)$ .

PROOF OF PROPOSITION 5.3 (i), (iii). Assertion (iii) follows from the fact that  $q_j(R)$  is less than the probability of getting to  $j$  with  $\xi$  starting at  $g+1$ . Combine this with the fact that

$$P_{\tau_j}^{(g+1, R)}(k_n \rightarrow G | k_n = j \text{ for some } n) \rightarrow 1$$

as  $j \rightarrow \infty$  to obtain (i).  $\square$

We must delay the proof of (ii) until we have proved the forthcoming Theorem 5.1.

For  $R \in \bigcup_{m=1}^{\infty} S_m$ , let

$$(5.8) \quad \begin{aligned} b(R) &= a(R) \quad \text{if } q_{a(R)+1}^{(R)} < q_{a(R)}^{(R)} \\ &= a(R) + 1 \quad \text{if } q_{a(R)+1}^{(R)} = q_{a(R)}^{(R)}. \end{aligned}$$

THEOREM 5.1. Suppose that  $G = \infty$ ,  $g > -\infty$ , and  $\lambda > .5$ . Let  $a(R) = b(R) = g+1$  for  $R \in S_{\infty}$ ,  $a(R) = b(R) = \infty$  for  $R \in S_0$ , and  $a$  and  $b$  be given by (5.6) and (5.8) on  $\bigcup_{m=1}^{\infty} S_m$ . A scheme  $\Psi$  is optimal for  $R_2$  if and only if

$$\begin{aligned} \Psi(k, R) &= g \quad \text{if } k < a(R) \\ &= R \quad \text{if } k \geq b(R). \end{aligned}$$

In particular,  $\Psi_2$ , given via (5.4) and (5.7), is optimal for  $R_2$ . The functions  $a$  and  $b$  are finite on  $S_0 \cup \bigcup_{m=1}^{\infty} S_m$ .

PROOF. With no loss of generality we take  $g = 0$ . Let

$$V_2(k, R) = P_{\Psi_2}^{(k, R)}(k_n \rightarrow \infty).$$

Clearly  $\Psi_2$  is optimal for  $S_{\infty}$  and, hence,

$$U(k; R, \lambda) = V_2(k, R) = E^R \left[ 1 - \left( \frac{1-\rho}{\rho} \right)^k \right] \quad \text{for } R \in S_{\infty} \text{ and all } k.$$

The reader may formally check or argue directly that  $f$  is not conserving, and therefore not optimal, at any  $(k, R)$  for which  $R \in S_{\infty}$ .

Clearly

$$V_2(k, R) = 1 - \left( \frac{1-\lambda}{\lambda} \right)^k \quad \text{for } R \in S_0 \text{ and all } k.$$

Condition (2.3) of Theorem 2.2 is obviously satisfied; the excessivity of  $V_2$ , and optimality of  $\Psi_2$ , for  $S_{\infty} \cup S_0$  is equivalent to

$$1 - \left( \frac{1-\lambda}{\lambda} \right)^k \geq E^R(\rho) E^{OR} \left[ 1 - \left( \frac{1-\rho}{\rho} \right)^{k+1} \right] + E^R(1-\rho) \left[ 1 - \left( \frac{1-\lambda}{\lambda} \right)^{k-1} \right]$$

for  $R \in S_0$  and all  $k$ , that is,

$$(5.9) \quad E^R(\rho) \leq \frac{2\lambda-1}{\lambda} \left( \frac{1-\lambda}{\lambda} \right)^{k-1} / E^{OR} \left[ \left( \frac{1-\lambda}{\lambda} \right)^{k-1} - \left( \frac{1-\rho}{\rho} \right)^{k+1} \right]$$

for  $R \in S_0$  and all  $k$ . We omit the straightforward calculation showing (5.9) to hold with strict inequality, and, therefore,  $f$  to be uniquely optimal, at each  $(k, R)$ ,  $R \in S_0$ .

Suppose, as an induction hypothesis, that  $\Psi_2$  is optimal for  $S_{\infty} \cup \bigcup_{i=0}^m S_i$ . For  $R \in S_{m+1}$ , we have, by Proposition 5.2, that  $\varphi R \in S_m$  and  $OR \in S_{\infty}$ . By the definition (5.6) of  $a(R)$  it is clear that  $f$ , and only  $f$ , is conserving at  $(k, R)$  for  $k < a(R)$  and that  $\varphi$  is conserving at  $(a(R), R)$ . By the proof of Theorem 3.3,  $R$  is conserving at  $(k, R)$  for  $k > a(R)$ . An appeal to Theorem 2.1 completes the proof by induction that  $\Psi_2$  is optimal.

Clearly  $g$  is conserving at  $(a(R), R)$  if  $b(R) = a(R) + 1$ . On the other hand, suppose  $g$  is conserving at  $(a(R), R)$ . Then, since  $R$  is conserving at  $(a(R) + 1, R)$  and Theorem 2.1 is applicable,  $\tau_{a(R)+1}$  is optimal. Hence  $b(R) = a(R) + 1$ .

Suppose, for a proof by contradiction which is very similar to the proof of Theorem 3.3, that  $\xi$  is optimal at some  $(k, R)$  for  $k > a(R)$ . Then an optimal strategy for that initial  $(k, R)$  is to pull  $g$  at stage 1, pull  $R$  at stage 2, and thereafter follow  $\Psi_2$ . An equally good, and therefore optimal, strategy is to pull  $R$  at stage 1, pull  $g$  at stage 2, and thereafter follow  $\Psi_2$ . But this latter strategy is not optimal since it can be improved by pulling  $R$  at stage 2 if  $Z_1 = 1$ .

That  $b(R) < \infty$  for  $R \in R_2 - S_0$  follows from Proposition 5.3. Even though the proof of Proposition 5.3 has not yet been completed, no circular reasoning is involved; for, only the optimality of  $\Psi_2$  from Theorem 5.1 will be needed for the completion of the proof of Proposition 5.3. ■

PROOF OF PROPOSITION 5.3 (ii). Since  $\Psi_2$  is optimal from Theorem 5.1,  $q_j(R)$  is at least as large as the probability of approaching  $\infty$  by pulling  $g$  until the fortune  $j$  is reached, then pulling  $R$  once, and thereafter pulling  $R$  or  $g$  accordingly as  $X_1 = 1$  or  $X_1 = -1$ ; that is (in case  $g = 0$ ),

$$q_j(R) \geq \frac{(2\lambda-1)/\lambda}{1-[(1-\lambda)/\lambda]^j} \{E^R(\rho)E^{OR} [1 - (\frac{1-\rho}{\rho})^{j+1}] + E^R(1-\rho)[1 - (\frac{1-\lambda}{\lambda})^{j-1}]\}$$

which, for  $R \in R_2 - S_0$ , can be shown to be larger than  $(2\lambda-1)/\lambda$  for sufficiently large  $j$ . ■

EXAMPLE 5.1. Suppose  $G = \infty$ ,  $g = 0$ ,  $\lambda = .6$ , and  $\rho$  is known to be either 0 or .625. The class

$$R'_2 = \{R: \text{supp } R \subset \{0, .625\}\}$$

is a subset of  $R_2$  and is closed under sampling. Let  $S'_m = S_m \cap R'_2$ ,  $m = 0, 1, \dots, \infty$ . Since  $(2\lambda-1)/\lambda = 1/3$ ,

$$S'_0 = \{R \in R'_2: R(\{.625\}) \leq 8/15\};$$

$$S'_m = \{R \in R'_2: \frac{8^m}{8^m + 7 \cdot 3^{m-1}} < R(\{.625\}) \leq \frac{8^{m+1}}{8^{m+1} + 7 \cdot 3^m}\}, m = 1, 2, \dots;$$

$$S'_\infty = \{R \in R'_2: R(\{.625\}) = 1\}.$$

Suppose  $R(\{.625\}) = 32/35$ . Then  $R \in S'_3$ ,  $\varphi R \in S'_2$ ,  $\varphi^2 R \in S'_1$ , and  $\varphi^3 R \in S'_0$ . Also  $\varphi R(\{.625\}) = 4/5$ ,  $\varphi^2 R(\{.625\}) = 3/5$ , and  $\varphi^3 R(\{.625\}) = 9/25$ . Straightforward calculations using (5.6) and (5.8) yield  $a(R) = b(R) = 1$ ,  $a(\varphi R) = b(\varphi R) = 3$ , and  $a(\varphi^2 R) = b(\varphi^2 R) = 14$ .

If the initial fortune is  $k = 2$ , then  $R$  should be pulled (since  $k \geq b(R)$ ). If  $X_1 = Z_1 = -1$ , so that  $k_1 = 1$  and  $R_1 = \varphi R$ , then  $R$  should be pulled until fortune 3 is reached (if ever), at which time  $R$  should be pulled again. If  $X_2 = -1$ , then  $R$  should be pulled until fortune 14 is reached, at which time  $R$  should be pulled again. If  $X_3 = -1$ , then  $R$  is never pulled again, no matter how large the fortune becomes. Of course,  $R$  is used exclusively if it ever gives a success.

If, instead, the initial fortune is  $k \geq 16$ , then  $R$  should be pulled the first three times. The modifications for other possible initial fortunes are now obvious.

In order to approach infinity from any starting point,  $\rho$  should always be pulled at least three times. Interestingly however, if  $k$  is small these pulls should be delayed until failures on  $\rho$  can be more easily withstood. ■

We turn to the case in which  $g$  and  $G$  are both finite. When  $G = \infty$ , the gambler, knowing that he wants to play forever, is quite interested in gathering information about  $R$ . Now that  $G < \infty$ , the gambler, especially when his fortune is close to  $G$ , may take a more immediate view. It is, therefore, not surprising to find that the gambler may now be wise to pull  $f$  when his fortune is close to  $G$ , thus eschewing the opportunity to learn about  $\rho$  that is provided by having a fortune far from  $g$ .

The definition (5.1) of  $S_{\infty}$  must be replaced by

$$(5.10) \quad S_{\infty} = \{R \in R_2 : E^R(\rho)E^{\sigma R}[u(k+1, \rho) - u(k-1, \lambda)] < u(k, \lambda) - u(k-1, \lambda) \text{ for all } k \in (g, G)\},$$

which is equivalent to (5.1) when  $G = \infty$ . We define  $S_m$  for  $m = 1, \dots, \infty$  via (5.2) and (5.3). Propositions 5.1 and 5.2 remain true for  $G < \infty$ . We continue to use (5.5) for  $g < j < G$ . We can also view  $\tau_j$  as used in (5.5) as a strategy for starting at  $(G-1, R)$ , as well as a strategy for starting at  $(g+1, R)$ . In analogy with (5.5), let

$$Q_j(R) = P_{\tau_j}^{(G-1, R)}(k_n \rightarrow G).$$

In analogy with (5.6), let

$$(5.11) \quad A(R) = \max\{k \in (g, G) : Q_k(R) = \max_{g < j < G} Q_j(R)\}.$$

Of course, in the present context, "max" and "g < j < G" appear in (5.6) also. Replace (5.7) by

$$(5.12) \quad \begin{aligned} \Psi_2(k, R) &= f \text{ if } k < a(R) \\ &= R \text{ if } a(R) \leq k \leq A(R) \\ &= g \text{ if } A(R) < k. \end{aligned}$$

For  $G < \infty$  there is no need for an analogue of Proposition 5.3. In analogy with (5.8), let

$$(5.13) \quad \begin{aligned} B(R) &= A(R) \text{ if } Q_{A(R)-1}(R) < Q_{A(R)}(R) \\ &= A(R)-1 \text{ if } Q_{A(R)-1}(R) = Q_{A(R)}(R) \end{aligned}$$

for  $R \in \bigcup_{m=1}^{\infty} S_m$ .

For  $G = \infty$  and any  $R \in R_2$  we saw that either  $f$  is uniquely optimal at  $(k, R)$  for all  $k$  or else there is some  $k$  such that  $f$  is not conserving at  $(k, R)$ . That this dichotomy does not hold when  $G < \infty$  is the reason for one further notation. Let

$$(5.14) \quad \begin{aligned} \bar{S}_0 &= \{R \in R_2 : E^R(\rho) E^{OR} [u(k+1, \rho) - u(k-1, \lambda)] \\ &\quad < u(k, \lambda) - u(k-1, \lambda) \text{ for all } k \in (g, G)\}. \end{aligned}$$

Clearly,  $S_0 \subset \bar{S}_0 \subset S_0 \cup S_1$ .

THEOREM 5.2. Suppose  $g$  and  $G$  are finite. Let  $a(R) = b(R) = g+1$  and  $A(R) = B(R) = G-1$  for  $R \in S_{\infty}$ ,  $a(R) = G$  and  $A(R) = g$  for  $R \in S_0$ , and  $a, b, A,$  and  $B$  be defined by (5.6), (5.8), (5.11), and

(5.13) on  $\bigcup_{m=1}^{\infty} S_m$ . A scheme  $\Psi$  is optimal for  $R_2$  if and only if

$$\begin{aligned} \Psi(k, R) &= f \text{ if } k < a(R) \\ &= R \text{ if } b(R) \leq k \leq B(R) \text{ and } R \in R_2 - \bar{S}_0 \\ &= g \text{ if } A(R) < k. \end{aligned}$$

In particular,  $\Psi_2$ , given via (5.4) and (5.12), is optimal for  $R_2$ .

For  $R \in R_2 - \bar{S}_0$ ,  $b(R) \leq B(R)$ .

We omit the proof which is very similar to the proof of Theorem 5.1.

EXAMPLE 5.2. Suppose  $g = 0$ ,  $G = 4$ ,  $\lambda = .74$ , and  $R(\{.75\}) = 1 - R(\{0\}) = .97554$ . Straightforward calculations show that  $R \in S_1 - \bar{S}_0$  and that  $a(R) = b(R) = A(R) = B(R) = 2$ . According to Theorem 5.2, in order for  $\Psi$  to be optimal it must satisfy  $\Psi(1,R) = \Psi(3,R) = g$  and  $\Psi(2,R) = \emptyset$ . This is a situation in which the set of fortunes where it is optimal to pull  $E$  is not an interval.  $\square$

Suppose that  $g = -\infty$ ,  $\lambda < .5$ , and  $R((\lambda, .5)) > 0$ . Then, according to the proof of Theorem 3.2,  $\{k: \Psi(k,R) = \emptyset\} \neq \emptyset$  for every optimal scheme. Therefore, there is no analogue of the set  $S_0$  and, as mentioned earlier, we pursue this case no further.

6. Optimal Schemes for the Class  $R_{\beta}$ . In this final section we consider the class

$$R_{\beta} = \{R: \text{supp } R \subset \{1-\beta, \beta\}\}$$

for the case  $G = \infty$ ,  $g > -\infty$ . The class  $R_{\beta}$  is closed under sampling and has the pleasant property that, for  $R \in R_{\beta}$ ,  $\sigma^s \varphi^f R$  depends on  $s$  and  $f$  only through  $s-f$ .

We carry over a notational convention from Example 1.1--if  $r = R(\{\beta\})$ , we not only use  $r$  as a number but we replace  $R$  by  $r$  in various symbols. For instance, we write  $U(k; r, \lambda)$  for  $U(k; R, \lambda)$  and  $\sigma^s \varphi^f r$  for  $(\sigma^s \varphi^f R)(\{\beta\})$ .

Since the case  $\beta = 1$  is discussed in Section 4 we take  $\beta \in (.5, 1)$  with no loss. Under the condition  $\lambda \leq .5$  or the condition  $\beta \leq \lambda$  the problem is trivial:  $R$  is optimal at every stage for the former and  $L$  is optimal at every stage for the latter. We therefore assume  $.5 < \lambda < \beta < 1$ . It seems reasonable to pull  $R$  if  $r$  is large enough and  $L$  otherwise. The next theorem says that some, but not all, optimal schemes have this characteristic. It indicates that, while a certain amount of flexibility is available for  $k$  large,  $R$  should never be pulled if  $r$  is small enough. The critical value of  $r$  is

$$r^* = \frac{(2\lambda-1)/\lambda}{(2\beta-1)/\beta} ,$$

which is the ratio  $u(g+1, \lambda)/u(g+1, \beta)$  where  $u$  is defined in Section 1.

**THEOREM 6.1.** Suppose  $G = \infty$ ,  $g > -\infty$ , and  $.5 < \lambda < \beta < 1$ .

Then a scheme  $\Psi$  is optimal for  $R_{\beta}$  if and only if

$$\begin{aligned} \Psi(k, r) &= R \text{ if } r > \left[ 1 + \left( \frac{1-\beta}{\beta} \right)^{k-g-1} \left( \frac{1-r^*}{r^*} \right)^{-1} \right] \\ &= L \text{ if } r < r^* . \end{aligned}$$

REMARK. The quantity

$$\left[1 + \left(\frac{1-\beta}{\beta}\right)^{k-g-1} \left(\frac{1-r^*}{r^*}\right)^{-1}\right]^{-1}$$

is equal to  $\sigma^{k-g-1} r^*$ ; it is the probability that  $\rho = \beta$  given  $k-g-1$  successes on  $R$  when the initial probability that  $\rho = \beta$  is  $r^*$ .

PROOF. With no loss of generality we may assume  $g = 0$ .

One of the schemes given by the theorem is

$$\begin{aligned} \Psi_{\beta}(k,r) &= R \quad \text{if } r \geq r^* \\ &= f \quad \text{if } r < r^* . \end{aligned}$$

Let

$$k^*(r) = \min\{j: \sigma^j r < r^*\} .$$

Since  $\sigma^s \varphi^f r$  depends only on  $s-f$  and  $r$  and, in view of the definition of  $k^*$ , the strategy induced by  $\Psi_{\beta}$  for starting at  $(k,r)$  specifies pulls of  $R$  until (if ever) reaching fortune  $0 \vee [k-k^*(r)]$ , and thereafter (if  $k > k^*(r)$ ) it specifies pulls of  $f$  indefinitely.

With the notation

$$V_{\beta}(k,r) = P_{\Psi_{\beta}}^{(k,r)}(k_n \rightarrow \infty)$$

we obtain

$$\begin{aligned} V_{\beta}(k,r) &= ru(k,\beta) \quad \text{if } k \leq k^*(r) \\ &= 1 - [1 - ru(k^*(r),\beta)][1 - u(k-k^*(r),\lambda)] \quad \text{if } k > k^*(r) . \end{aligned}$$

The remaining steps are: (i) to prove that  $V_{\beta}(k,r) \rightarrow 1$  as  $k \rightarrow \infty$  uniformly in  $r$ , that is, that (2.3) holds; (ii) to show that  $V_{\beta}$  is excessive for  $R_{\beta}$  and, thus, by Theorem 2.2, that  $\Psi_{\beta}$  is optimal for  $R_{\beta}$ ; (iii) to observe where equality holds in the calculations.

that show  $V_\beta$  to be excessive for  $R_\beta$  and, thus, conclude which schemes are conserving, and optimal by Theorem 2.1, for  $R_\beta$ . We shall carry out step (i) and leave the straightforward steps (ii) and (iii) to the reader.

Let  $\epsilon > 0$ . Choose  $K$  so that  $u(k, \lambda) > 1 - \epsilon$  for  $k \geq K/2$ ,  $u(k, \beta) > 1 - \epsilon/2$  for  $k \geq K/2$ , and  $r > 1 - \epsilon/2$  if  $k^*(r) \geq K/2$ . Suppose  $k \geq K$ . Either  $k - k^*(r) \geq K/2$  in which case  $1 - u(k - k^*(r), \lambda) < \epsilon$  or  $k^*(r) \geq K/2$  in which case  $1 - ru(k \wedge k^*(r), \beta) < \epsilon$ . In either case  $V_\beta(k, r) > 1 - \epsilon$ .  $\square$

### References

- Berry, Donald A. (1972). A Bernoulli two-armed bandit. Annals of Mathematical Statistics 43 871-897.
- Berry, Donald A. and Bert Fristedt (1979). Bernoulli one-armed bandits--arbitrary discount sequences. Annals of Statistics 7 1086-1105.
- Berry, Donald A. and Bert Fristedt (1980). Two-armed bandits with a goal; II: Dependent arms. Advances in Applied Probability 12 No. 4.
- Berry, Donald A., David C. Heath, and William D. Sudderth (1974). Red-and-black with unknown win probability. Annals of Statistics 2 602-608.
- Dubins, Lester E. and Leonard J. Savage (1976). Inequalities for Stochastic Processes: How to Gamble If You Must. Dover, New York.
- Dubins, Lester E. and William D. Sudderth (1977). Countably additive gambling and optimal stopping. Z. Wahrscheinlichkeitstheorie 41 59-72.
- Sudderth, William (1972). On the Dubins and Savage characterization of optimal strategies. Annals of Mathematical Statistics 43 498-507.