

Copyright © 1964, by the author(s).  
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

Electronics Research Laboratory  
University of California  
Berkeley, California  
Internal Technical Memorandum M-70

A NOTE ON COUNTER-EXAMPLES TO A  
CONJECTURE CONCERNING THE TWO-ARMED  
BANDIT

by

B. Gluss

\*This research reported herein is made possible through support received from the National Science Foundation grant G-15965/GP 2413.

May 19, 1964

A NOTE ON COUNTER-EXAMPLES TO A CONJECTURE  
CONCERNING THE TWO-ARMED BANDIT\*

Brian Gluss

University of California, Berkeley

ABSTRACT

In the case of a two-armed bandit comprising two Bernoulli machines with known parameters  $p_1$  and  $p_2$ , where it is not known which parameter pertains to which machine, Feldman has proved that, to maximize the total expected return over  $r$  trials, that machine with the higher a posteriori parameter should be chosen at each trial. It has been conjectured that an analogous result holds when only probability distributions for each machine's parameter are known, selecting that machine with the higher expected a posteriori parameters. It is shown in this note that this conjecture is not generally correct.

---

\*The research reported herein is made possible through support received from the National Science Foundation grant G-15965/GP 2413.

A NOTE ON COUNTER-EXAMPLES TO A CONJECTURE  
CONCERNING THE TWO-ARMED BANDIT\*

Brian Gluss

Electronics Research Laboratory  
University of California, Berkeley

INTRODUCTION

Consider a two-armed bandit comprising two Bernoulli machines with known parameters  $p_1$  and  $p_2$ , where it is not known which parameter pertains to which machine. Feldman<sup>1</sup> has proved that in order to maximize the total expected return over  $r$  trials,  $r$  fixed, in which a trial produces scores

$$\left. \begin{array}{ll} 1 & \text{with probability } p_i \\ 0 & \text{with probability } 1 - p_i \end{array} \right\} \quad (1)$$

according as to which machine  $i$  ( $= 1, 2$ ) is used, the optimal policy is to select at each trial that machine which has the higher a posteriori probability  $p$  of scoring a 1.

Suppose instead that, as in Gluss,<sup>2</sup> the Bernoulli machines I and II, have a priori distributions for their parameters  $p$  and  $p'$  respectively of the forms

---

\*The research reported herein is made possible through support received from the National Science Foundation grant G-15965/GP 2413.

$$dG(p) = \frac{p^{a-1}(1-p)^{b-1}}{B(a, b)} dp, \quad (2)$$

and

$$dG'(p) = \frac{(p')^{a'-1}(1-p')^{b'-1}}{B(a', b')} dp', \quad (3)$$

so that, after  $m$  1's and  $n$  0's with  $I$ , the Bayes a posteriori distribution for  $p$  is given by

$$dG_{m, n}(p) = \frac{p^m(1-p)^n dG(p)}{\int_0^1 x^m(1-x)^n dG(x)},$$

i. e. ,

$$dG_{m, n}(p) = \frac{p^{a+m-1}(1-p)^{b+n-1}}{B(a+m, b+n)} dp, \quad (4)$$

with expectation

$$p_{m, n} = \frac{m+a}{(m+a) + (n+b)}, \quad (5)$$

with similar expressions for  $dG'_{m', n'}(p')$  and  $p'_{m', n'}$ . The  $B(u, v)$  are beta distributions.

In this case, it has been conjectured by L. A. Zadeh and others that a result analogous to that in Feldman's model holds. That is, in order to maximize the total expected return for an  $r$ -stage process, the optimal policy is to choose at each trial that machine with the higher expected a posteriori probability of scoring a 1. That is,

$$\text{choose machine I if and only if } \frac{m+a}{(m+a) + (n+b)} > \frac{m'+a'}{(m'+a') + (n'+b')}. \quad (6)$$

We shall show below that this conjecture does not always hold, although it does over large regions of the  $(m, n, m', n')$  space.

### THE FUNCTIONAL EQUATION

It will be notationally convenient to introduce

$$\left. \begin{aligned} M = m+a, & \quad N = n+b, & R = M + N, \\ M' = m'+a', & \quad N' = n'+b', & R' = M' + N'. \end{aligned} \right\} \quad (7)$$

The policy of Eq. (6) then becomes

$$\text{choose I if and only if} \quad \frac{M}{R} > \frac{M'}{R'} \quad (8)$$

Let  $f_{M, N, M', N'}^r$  = expected total return for a process with  $r$  trials remaining, after  $M - a$  1's and  $N - b$  0's with I and  $M' - a'$  1's and  $N' - b'$  0's with II, using an optimal policy.

Then we have the functional equation

$$f_{M, N, M', N'}^r = \text{Max} \left[ \begin{aligned} \text{I:} & \int_0^1 p \left[ 1 + f_{M+1, N, M', N'}^{r-1} \right] dG_{m, n}(p) \\ & + \int_0^1 (1-p) f_{M, N+1, M', N'}^{r-1} dG_{m, n}(p); \\ \text{II:} & \int_0^1 p' \left[ 1 + f_{M, N, M'+1, N'}^{r-1} \right] dG'_{m', n'}(p') \\ & + \int_0^1 (1-p') f_{M, N, M', N'+1}^{r-1} dG'_{m', n'}(p') \end{aligned} \right]$$

That is,

$$f_{M, N, M', N'}^r = \text{Max} \left[ \begin{array}{l} \text{I: } \frac{M}{R} \left[ 1 + f_{M+1, N, M', N'}^{r-1} \right] + \frac{N}{R} f_{M, N+1, M', N'}^{r-1} \\ \text{II: } \frac{M'}{R'} \left[ 1 + f_{M, N, M'+1, N'}^{r-1} \right] + \frac{N'}{R'} f_{M, N, M', N'+1}^{r-1} \end{array} \right] \quad (9)$$

We have the boundary conditions  $f^0 = 0$ , and the policy of Eq. (8) is evidently optimal for  $r = 1$ , and

$$f_{M, N, M', N'}^1 = \text{Max} \left[ \frac{M}{R}, \frac{M'}{R'} \right]. \quad (10)$$

#### COUNTER-EXAMPLES FOR $r = 2$

We now consider the two-stage process ( $r = 2$ ), and assume that initially

$$\frac{M}{R} > \frac{M'}{R'}. \quad (11)$$

It will be shown that the policy of Eq. (8) at the first trial--i. e., choose machine I--is not optimal for all  $M, M', N, N'$ .

Using Eq. (10), and remembering that inequality (11) implies that we also have

$$\text{and} \left. \begin{array}{l} \frac{M+1}{R+1} > \frac{M'}{R'} \\ \frac{M}{R} > \frac{M'}{R'+1} \end{array} \right\} \quad (12)$$

Eq. (9) reduces to

$$f^2_{M, N, M', N'} = \text{Max} \left[ \begin{array}{l} \text{I: } \frac{M}{R} \left[ 1 + \frac{M+1}{R+1} \right] + \frac{N}{R} \text{Max} \left\{ \frac{M}{R+1}, \frac{M'}{R'} \right\}; \\ \text{II: } \frac{M'}{R'} \left[ 1 + \text{Max} \left\{ \frac{M}{R}, \frac{M'+1}{R'+1} \right\} \right] + \frac{N'}{R'} \left[ \frac{M}{R} \right] \end{array} \right] \quad (13)$$

We discuss below the four cases  $\frac{M}{R+1} < \frac{M'}{R'}$ ;  $\frac{M}{R} < \frac{M'+1}{R'+1}$ , all consistent with inequality (11), determining expressions for I - II, and showing that if

$$\frac{M}{R} \geq \frac{M'+1}{R'+1}, \quad (14)$$

policy (8) is optimal, while otherwise this is not necessarily the case. That is, for

$$\frac{M'}{R'} < \frac{M}{R} < \frac{M'+1}{R'+1}, \quad (15)$$

policy (8) is not necessarily optimal.

$$(1) \quad \underline{\frac{M}{R} \geq \frac{M'+1}{R'+1}}, \quad \frac{M}{R+1} < \frac{M'}{R'}$$

In this case,

$$\begin{aligned} \text{I} - \text{II} &= \frac{M}{R} + \frac{M(M+1)}{R(R+1)} + \frac{N}{R} \cdot \frac{M'}{R'} - \frac{M'}{R'} - \frac{M'}{R'} \cdot \frac{M}{R} - \frac{N'}{R'} \cdot \frac{M}{R} \\ &= \frac{M}{R} - \frac{M'}{R'} + \frac{M(M+1)}{R(R+1)} + \frac{N}{R} \cdot \frac{M'}{R'} - \frac{M}{R} \\ &= \frac{M}{R} - \frac{M'}{R'} - \frac{M}{R} \cdot \frac{N}{R+1} + \frac{N}{R} \cdot \frac{M'}{R'} \\ &= \left( \frac{M}{R} - \frac{M'}{R'} \right) + \frac{N}{R} \left( \frac{M'}{R'} - \frac{M}{R+1} \right). \end{aligned}$$



Since  $\frac{M'}{R'} > \frac{M}{R+1}$  and  $\frac{M}{R} > \frac{M'}{R'}$ , here we have I > II.

$$(2) \quad \frac{M}{R} \geq \frac{M'+1}{R'+1}, \quad \frac{M}{R+1} > \frac{M'}{R'}$$

We now have

$$\begin{aligned} I - II &= \frac{M}{R} + \frac{M(M+1)}{R(R+1)} + \frac{N}{R} \cdot \frac{M}{R+1} - \frac{M'}{R'} - \frac{M'}{R'} \cdot \frac{M}{R} - \frac{N'}{R'} \cdot \frac{M}{R} \\ &= \frac{M}{R} - \frac{M'}{R'}, \end{aligned}$$

which is once again greater than zero. We now reach the two questionable cases.

$$(3) \quad \frac{M}{R} < \frac{M'+1}{R'+1}, \quad \frac{M}{R+1} < \frac{M'}{R'}$$

Here we have

$$\begin{aligned} I - II &= \frac{M}{R} + \frac{M(M+1)}{R(R+1)} + \frac{N}{R} \cdot \frac{M'}{R'} - \frac{M'}{R'} - \frac{M'(M'+1)}{R'(R'+1)} - \frac{N'}{R'} \cdot \frac{M}{R} \\ &= \frac{M(M+1)}{R(R+1)} - \frac{M'(M'+1)}{R'(R'+1)}. \end{aligned}$$

Hence, if

$$\frac{M(M+1)}{R(R+1)} < \frac{M'(M'+1)}{R'(R'+1)}, \quad \left. \begin{array}{l} \frac{M}{R} < \frac{M'+1}{R'+1}, \\ \frac{M}{R+1} < \frac{M'}{R'}, \end{array} \right\} \quad (16)$$

$$\frac{M}{R} < \frac{M'+1}{R'+1},$$

$$\frac{M}{R+1} < \frac{M'}{R'},$$

and

$$\frac{M}{R} > \frac{M'}{R'},$$

once again the conjecture is false. We have thus found that, while it holds under all other conditions, the conjecture fails under the two sets of conditions (16) and (17).

and

$$\left. \begin{aligned} \frac{M}{M'} &> \frac{R}{R'} \\ \frac{M}{M'} &> \frac{R+1}{R'} \\ \frac{M}{M'} &> \frac{R}{R'+1} \\ \frac{M}{M'} - \frac{R}{R'} &> \frac{R}{M'} \left( \frac{R+1}{M'+1} - \frac{R}{M} \right) \end{aligned} \right\} (17)$$

When this expression is greater than zero, the conjecture holds. However, when

$$I - II = \frac{R}{M} + \frac{R(R+1)}{M(M+1)} + \frac{R}{N} \cdot \frac{R+1}{M} - \frac{R}{M'} - \frac{R'(R'+1)}{M'(M'+1)} - \frac{R}{N} \cdot \frac{R}{M} = \left( \frac{R}{M} - \frac{R}{M'} \right) + \frac{R}{M'} \left( \frac{R}{M} - \frac{R}{M'+1} \right)$$

In this final case,

$$(4) \quad \frac{R}{M} > \frac{R}{M'+1}, \quad \frac{R+1}{M} > \frac{R}{R'}$$

$$\frac{M(M+1)}{M'(M'+1)} > \frac{R(R+1)}{R'(R'+1)}$$

which are mutually consistent inequalities, in this situation the conjecture does not hold. Of course, in case (3), the conjecture holds when

## TWO NUMERICAL EXAMPLES

(a) For  $M = 17$ ,  $N = 11$ ,  $M' = 3$ ,  $N' = 2$ , the inequalities (16) hold, since

$$\frac{17.18}{28.29} < \frac{3.4}{5.6}, \quad \frac{17}{28} < \frac{4}{6}, \quad \frac{17}{29} < \frac{3}{5}, \quad \text{and} \quad \frac{17}{28} > \frac{3}{5}.$$

In this case

$$I = \frac{17}{28} \left(1 + \frac{18}{29}\right) + \frac{11}{28} \cdot \frac{3}{5} = \frac{1}{70} \left(85 \frac{11}{29}\right).$$

$$II = \frac{3}{5} \left(1 + \frac{4}{6}\right) + \frac{3}{5} \cdot \frac{17}{28} = \frac{1}{70} (87),$$

so that although  $\frac{17}{28} > \frac{3}{5}$ ,  $I < II$ .

(b) Let  $M = 50$ ,  $N = 31$ ,  $M' = 3$ ,  $N' = 2$ . Then the inequalities (17) hold, since

$$\frac{50}{81} - \frac{3}{5} < \frac{3}{5} \left(\frac{4}{6} - \frac{50}{81}\right), \quad \frac{50}{81} < \frac{4}{6}, \quad \frac{50}{82} > \frac{3}{5}, \quad \text{and} \quad \frac{50}{81} > \frac{3}{5}.$$

Here we have

$$I = \frac{50}{81} \left(1 + \frac{51}{82}\right) + \frac{31}{81} \cdot \frac{50}{82} = \frac{100}{81}.$$

$$II = \frac{3}{5} \left(1 + \frac{4}{6}\right) + \frac{2}{5} \cdot \frac{50}{81} = \frac{101}{81}.$$

## GENERAL COMMENTS

For general  $r$ , one might hope for simple conditions under which the maximum-expected-a-posteriori-policy holds, and correspondingly simple conditions under which it does not. However, as  $r$  increases

these conditions become progressively more complicated, comprising larger and larger numbers of inequalities on  $M, R, M', R'$ . Moreover, the sum of these regions does not appear to tend to a null region as  $r \rightarrow \infty$  (when we introduce a discount factor into the model), so that even in the asymptotic case of an infinite number of trials, the conjecture appears false.

It might be of interest to give--for brevity, without proofs--some simple regions for which the conjecture holds.

$$(i) \quad \frac{M}{R+r-1} > \frac{M'+r-1}{R'+r-1} \quad (18)$$

This condition is intuitively obvious, since the a posteriori probabilities of machine I are then always greater than those of machine II, whichever machine is used at any trial and whatever the result of its use.

$$(ii) \quad \frac{M}{R+r-1} > \frac{M'}{R'} \quad \text{and} \quad N+r-2 < M'. \quad (19)$$

$$(iii) \quad \frac{M}{R} > \frac{M'+r-1}{R'+r-1} \quad \text{and} \quad N > M'+r-2. \quad (20)$$

In both these cases, (ii) and (iii), the inequalities imply that

$$\frac{M}{R+r-1-i} > \frac{M'+i}{R'+i}, \quad i = 0, 1, \dots, r-1, \quad (21)$$

so that losing at each stage with machine I still leaves us with higher a posteriori probabilities than winning at each stage with machine II.

## REFERENCES

1. FELDMAN, D., (1962). Contributions to the two-armed bandit problem. Annals of Math. Stat. 33, 847-856.
2. GLUSS, B., (1962). A note on a computational approximation to the two-machine problem. Inf. and Control 5, 268-275.