

Finite-Time Lower Bounds for the Two-Armed Bandit Problem

Sanjeev R. Kulkarni and Gábor Lugosi

Abstract—We obtain minimax lower bounds on the regret for the classical two-armed bandit problem. We provide a finite-sample minimax version of the well-known $\log n$ asymptotic lower bound of Lai and Robbins. The finite-time lower bound allows us to derive conditions for the amount of time necessary to make any significant gain over a random guessing strategy. These bounds depend on the class of possible distributions of the rewards associated with the arms. For example, in contrast to the $\log n$ asymptotic results on the regret, we show that the minimax regret is achieved by mere random guessing under fairly mild conditions on the set of allowable configurations of the two arms. That is, we show that for every allocation rule and for every n , there is a configuration such that the regret at time n is at least $1 - \epsilon$ times the regret of random guessing, where ϵ is any small positive constant.

Index Terms—Bandit, estimation, learning, lower bound, minimax, regret, two-armed.

I. INTRODUCTION

Bandit problems have received considerable interest (e.g., see [5] and [6]) originating from the work in [8] and [9]. In the classical two-armed bandit problem, there are two unknown distributions P_1 and P_2 associated with arm 1 and arm 2, respectively. At each time we are allowed to select an arm from which to receive a reward drawn independently according to the distribution for that arm. Our goal is to maximize the expected sum of the rewards. Let m_1 and m_2 denote the expected values corresponding to P_1 and P_2 , respectively. If we knew which one of m_1 or m_2 is larger, we could keep selecting the arm with larger mean, and after time n , our expected reward would be $n \max(m_1, m_2)$. Since the distributions P_1 and P_2 are unknown, the expected reward will always be smaller than this optimal value. The difference between $n \max(m_1, m_2)$ and the expected reward is called the regret. Note that if, in each step, we select an arm independently with equal probabilities, the regret is $n\Delta/2$, where $\Delta = |m_1 - m_2|$. The results of Lai and Robbins [7] and subsequent extensions by others (e.g., [1]–[4]) showed that in a fairly strong asymptotic sense the optimum achievable regret is $\Delta \log n / \bar{I}$, where \bar{I} is the Kullback–Leibler divergence between P_1 and P_2 .

In this paper, we consider the problem from a nonasymptotic minimax perspective. In the minimax setup, one assumes that the possible configurations (P_1, P_2) belong to a fixed family of pairs of distributions. Then one is interested in the smallest possible regret achievable by any strategy of selecting the arms. Here we obtain lower bounds on the regret of such optimal strategies. Our starting point is a class of two configurations, and we show that even if the player knows this class, he/she cannot perform better than the lower bounds obtained here.

First we offer a finite-sample minimax version of the Lai–Robbins lower bound (see Theorem 1 below). This result can be used to provide bounds on the sample size necessary to guarantee a desired performance. Also, in sharp contrast to the well-known $\log n$ asymptotic

results on the regret, we show that for small sample sizes the wrong arm will be pulled a constant fraction of times. In particular, the minimax regret is about $n\Delta/2$ under fairly mild conditions on the set of allowable configurations of the two arms. We show that if the set of allowable configurations is sufficiently “large,” then for any n , for any small ϵ , and for any strategy of selecting arms, there is a configuration such that the regret is larger than $(1 - \epsilon)n\Delta/2$. In other words, regardless of how large n is, up to time n , the “bad” arm will be played almost half of the time for some configuration. That is, in the minimax sense, no arm-selection strategy can perform better than completely random selections.

II. FORMULATION AND LOWER BOUNDS

Let a configuration $\Theta = (\theta_1, \theta_2)$ be a pair of parameters determining the distributions of the two arms. That is, if arm i is pulled, a reward is paid independently according to the probability density f_{θ_i} ($i = 1, 2$). All densities are understood with respect to a common dominating σ -finite measure λ on the real line. Denote the respective means by

$$m_i = \int x f_{\theta_i}(x) d\lambda(x), \quad i = 1, 2$$

and let $\Delta = |m_1 - m_2|$. Assume without loss of generality that $m_1 > m_2$, that is, arm 1 is optimal. We denote the measure and expectation with respect to the distribution associated with θ by P_θ and E_θ , respectively.

Introduce the alternative configuration $\Theta' = (\theta'_1, \theta_2)$, for some θ'_1 such that $m'_1 = \int x f_{\theta'_1}(x) d\lambda(x) = m_2 - \Delta$. Thus, the distribution of the reward after pulling arm 2 is unchanged, but arm 2 is optimal in configuration Θ' .

Our bounds involve the *information divergence* (or Kullback–Leibler number) between the densities f_{θ_1} and $f_{\theta'_1}$ given by

$$\begin{aligned} I &= I(\theta'_1, \theta_1) = \int f_{\theta'_1}(x) \log \left(\frac{f_{\theta'_1}(x)}{f_{\theta_1}(x)} \right) d\lambda(x) \\ &= E_{\theta'_1} \left\{ \log \left(\frac{f_{\theta'_1}(x)}{f_{\theta_1}(x)} \right) \right\} \end{aligned}$$

as well as a variance-like quantity, denoted V , related to the information divergence I by

$$\begin{aligned} V &= \int f_{\theta'_1}(x) \log^2 \left(\frac{f_{\theta'_1}(x)}{f_{\theta_1}(x)} \right) d\lambda(x) - I^2 \\ &= E_{\theta'_1} \left\{ \log^2 \left(\frac{f_{\theta'_1}(x)}{f_{\theta_1}(x)} \right) \right\} - I^2. \end{aligned}$$

(All logarithms are of the natural base.)

An adaptive allocation rule $\Phi = (\phi_1, \phi_2, \dots)$ is a sequence of random variables taking values in $\{1, 2\}$ such that ϕ_i is measurable with respect to the σ -field \mathcal{F}_{i-1} generated by the previous values $\phi_1, X_1, \dots, \phi_{i-1}, X_{i-1}$, where X_1, X_2, \dots are the random variables denoting the sequence of rewards obtained. That is, based on the previous rewards (X_1, \dots, X_{i-1}) and the previous selections $(\phi_1, \dots, \phi_{i-1})$, ϕ_i denotes whether arm 1 or arm 2 is to be pulled at time i . Under a particular adaptive allocation rule and configuration Θ , our reward up to and including time n is

$$S_n = \sum_{i=1}^n X_i.$$

Manuscript received November 5, 1997, revised June 1, 1999. Recommended by Associate Editor, D. Yao. This work was supported in part by the National Science Foundation under NYI Grant IRI-9457645. The work of G. Lugosi was supported by DIGES under Grant PB96-0300.

S. R. Kulkarni is with the Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 USA (e-mail: kulkarni@ee.princeton.edu).

G. Lugosi is with the Department of Economics, Pompeu Fabra University, Ramon Trias Fargas 25-27, 08005 Barcelona, Spain (e-mail: lugosi@upf.es).

Publisher Item Identifier S 0018-9286(00)04082-4.

Since $E[X_i | \mathcal{F}_{i-1}] = m_{\phi_i}$, the expected reward is

$$E[S_n] = \sum_{i=1}^n E[X_i] = \sum_{i=1}^n E(E[X_i | \mathcal{F}_{i-1}]) = E \left[\sum_{i=1}^n m_{\phi_i} \right]$$

and the regret at time n is

$$R_n(\Theta) = n \max(m_1, m_2) - E[S_n].$$

In other words, $R_n(\Theta)$ is Δ times the expected number of times the arm with worse expected payoff is pulled.

Our goal is to obtain lower bounds on the minimum value of

$$\max(R_n(\Theta), R_n(\Theta'))$$

over all possible adaptive allocation rules. For any integer $n > 0$ and $\alpha > 0$, introduce

$$\Lambda_{\alpha, n} = e^{-\alpha I} \left(1 - \max_{1 \leq i \leq n} p_{\alpha, i} \right)$$

where

$$p_{\alpha, n} = P_{\theta'} \left\{ \sum_{i=1}^n \log \left(\frac{f_{\theta_1}(x_i)}{f_{\theta'_1}(x_i)} \right) \leq -\alpha I \right\}.$$

Theorem 1: For any $n, a_n \in (0, 1), c_n \in (0, n), \alpha > c_n$, and for any adaptive allocation rule

$$\begin{aligned} \max(R_n(\Theta), R_n(\Theta')) \\ \geq \Delta \min(c_n(1 - a_n), (n - c_n)a_n \Lambda_{\alpha, c_n}). \end{aligned}$$

To interpret the theorem, we need lower bounds for $\Lambda_{\alpha, n}$, that is, upper bounds for $p_{\alpha, n}$. Note that $p_{\alpha, n}$ is the probability that the sum of n i.i.d. random variables (with negative mean $-I$) is less than the mean of the sum minus $(\alpha - n)I$. Thus, it follows by Chebyshev's inequality that

$$p_{\alpha, i} \leq \frac{iV}{(\alpha - i)^2 I^2}$$

and therefore, since the right-hand side is a monotonically increasing function of i , we have

$$\Lambda_{\alpha, n} \geq e^{-\alpha I} \left(1 - \frac{nV}{(\alpha - n)^2 I^2} \right). \quad (1)$$

We will see that (1) is a satisfactory bound if the ratio V/I is not too large. This is indeed the case for many interesting cases. The next two examples serve as illustration.

Example: Let f_{θ_1} be the normal density with mean m_1 and variance σ^2 , and let $f_{\theta'_1}$ be the normal density with mean m'_1 and variance σ^2 . Then straightforward calculation shows that $V = I$ for all values of m_1, m'_1 , and σ .

Example: Let θ_1 correspond to the Bernoulli distribution $P_{\theta_1}(\{0\}) = p, P_{\theta_1}(\{1\}) = 1 - p$, and let θ'_1 be defined by $P_{\theta'_1}(\{0\}) = 1 - p, P_{\theta'_1}(\{1\}) = p$ and assume that $p > 1/2$. Then using the inequality $\log x \leq x - 1$

$$\frac{V}{I} = \frac{(1 - (2p - 1)^2) \log^2(p/(1 - p))}{(2p - 1) \log(p/(1 - p))} \leq 4p \leq 4.$$

Note that in this example we can take f_{θ_2} to be any density with mean $m_2 = 1/2$, for example, we may let $P_{\theta_2}(\{1/2\}) = 1$.

In specific situations, one may get much sharper estimates. For example, if both f_{θ_1} and $f_{\theta'_1}$ are Gaussian with variance σ^2 , then $\log(f_{\theta_1}(X)/f_{\theta'_1}(X))$ also has a Gaussian distribution, so one may get sharper estimates for $p_{\alpha, n}$ by using standard bounds for the

tail of a Gaussian distribution, but we do not detail these, rather straightforward, bounds here.

Corollary 1: Fix any $\epsilon \in (0, 1)$. If n is so large that

$$n\epsilon^2 \geq \max \left(4 \frac{(1 - \epsilon)^2 \log n}{\epsilon I}, e^{2V/(I(1 - \epsilon))} \right)$$

then

$$\max(R_n(\Theta), R_n(\Theta')) \geq \Delta \left[(1 - \epsilon)^2 \frac{\log n}{I} \right].$$

Proof: In Theorem 1 take $a_n = \epsilon, c_n = \lfloor (1 - \epsilon) \log n / I \rfloor$, and $\alpha = (1 + \epsilon)c_n$. Then (1) and a straightforward calculation shows that

$$c_n(1 - a_n) \leq (n - c_n)a_n \Lambda_{\alpha, c_n}$$

whenever the condition for n is satisfied, and therefore $\max(R_n(\Theta), R_n(\Theta')) \geq c_n(1 - a_n)$. \square

Remark: The classical lower bound of Lai and Robbins [7] states that for any allocation rule, the regret $R_n(\Theta)$ is asymptotically not smaller than $(\Delta/\tilde{I}) \log n$, where \tilde{I} is the Kullback–Leibler divergence between P_1 and P_2 . The main novelty of Corollary 1 is in its nonasymptotic nature. Note however, that the new bound is for $\max(R_n(\Theta), R_n(\Theta'))$ instead of just $R_n(\Theta)$, and the divergence I appearing in the bound is different from \tilde{I} . However, it is easy to see that by adding extra regularity conditions similar to those of Lai and Robbins, the bound of Corollary 1 may easily be converted into a nonasymptotic analogue of the Lai–Robbins lower bound.

For smaller values of n , Theorem 1 may be used to derive much larger lower bounds.

Corollary 2: Let $a \in (0, 1)$. If $n^{1-a} \geq 4e$ and $n^a I + \sqrt{n^a V/2} \leq 1$ then

$$\max(R_n(\Theta), R_n(\Theta')) \geq \Delta \frac{\lfloor n^a \rfloor}{2}.$$

Proof: Take $c_n = \lfloor n^a \rfloor, a_n = 1/2$, and $\alpha = n + \sqrt{nV/(2I^2)}$ in Theorem 1. \square

Taking $c_n = n/2$ and $a_n = 1/2$ in Theorem 1 we obtain the lower bound $(n\Delta/4)\Lambda_{\alpha, n/2}$. The following theorem improves on this.

Theorem 2: For any $n, \alpha > n$, and for any adaptive allocation rule

$$\max(R_n(\Theta), R_n(\Theta')) \geq \frac{n\Delta}{2} \Lambda_{\alpha, n}.$$

The next corollary points out that for sample sizes of the order of $1/I$, the number of times the bad armed is pulled is linear. This means that for sample sizes smaller than this, it is basically impossible to learn which is the best arm.

Corollary 3: Let c be a positive constant. For any $n \leq c/I$ and for any adaptive allocation rule

$$\max(R_n(\Theta), R_n(\Theta')) \geq \frac{\Delta n}{4} e^{-c - \sqrt{cV/(2I)}}.$$

Proof: Note that if we take $\alpha = n + \sqrt{nV/(2I^2)}$

$$1 - \frac{nV}{(\alpha - n)^2 I^2} = \frac{1}{2}$$

and therefore, the corollary follows by applying Theorem 2 with (1) for $n \leq c/I$. \square

Corollary 4: Let $\epsilon > 0$ be arbitrary. If there exists an $\alpha > n$ such that $\alpha \leq -\log(1 - \epsilon)/2I$ and $(\alpha - n)^2/n \geq 2V/(\epsilon I^2)$, then

$$\max(R_n(\Theta), R_n(\Theta')) \geq \frac{\Delta n}{2} (1 - \epsilon).$$

Proof: Straightforward calculation shows that if the conditions are satisfied then $e^{-\alpha I} \geq \sqrt{1-\epsilon}$ and $1 - nV / ((n-\alpha)^2 I^2) \geq \sqrt{1-\epsilon}$, so the statement follows by Theorem 2 and the bounds of (1). \square

Corollary 4 may be applied with arbitrary ϵ in many cases when, in the class of allowable configurations, there are pairs (θ_1, θ'_1) with arbitrarily small information divergence. The following two special cases illustrate such situations.

Corollary 5: Suppose P_2 is an arbitrary distribution with mean zero (which can even be known). Suppose P_1 is Gaussian with mean either Δ or $-\Delta$ with arbitrary variance. Then for every adaptive allocation rule, for every $\epsilon > 0$, and for every n , there is a configuration such that the regret at time n is at least $(1-\epsilon)n\Delta/2$.

Corollary 6: Let S be the set of all configurations $\Theta = (\theta_1, \theta_2)$ such that both P_{θ_1} and P_{θ_2} are Bernoulli distributions (i.e., of the form $P_{\theta}(\{0\}) = p$, $P_{\theta}(\{1\}) = 1-p$ for some p). Then for every adaptive allocation rule Φ , for every $\epsilon > 0$, and for every n

$$\sup_{\Theta \in S} R_n(\Theta) \geq \frac{n\Delta}{2}(1-\epsilon)$$

where Δ is the difference between the means corresponding to the two arms.

III. A CHANGE-OF-MEASURE LEMMA

As before, let the vector $\mathbf{X} = (X_1, \dots, X_n)$ of random variables denote the rewards up to time n if a particular adaptive allocation rule is used. Let $\mathbf{x} = (x_1, x_2, \dots, x_n) \in \mathcal{R}^n$ denote a fixed realization of \mathbf{X} .

$T_{\mathbf{x}}(1)$ and $T_{\mathbf{x}}(2)$ denote the number of times arm 1, and arm 2 are pulled up to time n .

The key part of the proofs of the results in the previous section is the following measure-transformation lemma, which is based on ideas of Lai and Robbins [7].

Lemma 1: For any integer $k \in [0, n]$, and $\alpha > 0$

$$P_{\Theta}\{T_{\mathbf{x}}(1) = k\} \geq e^{-\alpha I}(1-p_{\alpha,k})P_{\Theta'}\{T_{\mathbf{x}}(1) = k\}.$$

Proof: Let $J \subset \{1, \dots, n\}$ be a set of indices. On J , introduce the likelihood ratio

$$L_J(\mathbf{x}) = \sum_{j \in J} \log \left(\frac{f_{\theta_1}(x_j)}{f_{\theta'_1}(x_j)} \right).$$

The first step of the proof is trivial

$$P_{\Theta}\{T_{\mathbf{x}}(1) = k\} \geq P_{\Theta}\{T_{\mathbf{x}}(1) = k, L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I\} \quad (2)$$

where $B(\mathbf{x})$ is the set of indices indicating the times when arm 1 is pulled by the allocation rule based on the sequence of observations \mathbf{x} .

For each index set J , define $A_J \subset \mathcal{R}^n$ by

$$A_J = \{\mathbf{x} : B(\mathbf{x}) = J, L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I\}.$$

Thus

$$\begin{aligned} P_{\Theta}\{T_{\mathbf{x}}(1) = k, L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I\} \\ = P_{\Theta} \left\{ \bigcup_{J:|J|=k} A_J \right\} = \sum_{J:|J|=k} P_{\Theta}\{A_J\}. \end{aligned}$$

Now

$$\begin{aligned} P_{\Theta}\{A_J\} &= \int_{A_J} \left(\prod_{j \in J} f_{\theta_1}(x_j) \right) \left(\prod_{j \notin J} f_{\theta_2}(x_j) \right) \\ &\quad \times d\lambda(x_1) \dots d\lambda(x_n) \\ &= \int_{A_J} \left(\prod_{j \in J} \frac{f_{\theta_1}(x_j)}{f_{\theta'_1}(x_j)} \right) \left(\prod_{j \in J} f_{\theta'_1}(x_j) \right) \\ &\quad \times \left(\prod_{j \notin J} f_{\theta_2}(x_j) \right) d\lambda(x_1) \dots d\lambda(x_n). \end{aligned}$$

But for each $\mathbf{x} \in A_J$, we have $L_J(\mathbf{x}) > -\alpha I$, so

$$\left(\prod_{j \in J} \frac{f_{\theta_1}(x_j)}{f_{\theta'_1}(x_j)} \right) > e^{-\alpha I}$$

and therefore

$$P_{\Theta}\{A_J\} \geq e^{-\alpha I} P_{\Theta'}\{A_J\}.$$

It follows that

$$\begin{aligned} P_{\Theta}\{T_{\mathbf{x}}(1) = k, L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I\} \\ \geq e^{-\alpha I} P_{\Theta'}\{T_{\mathbf{x}}(1) = k, L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I\} \\ = e^{-\alpha I} P_{\Theta'}\{T_{\mathbf{x}}(1) = k\} \\ \times P_{\Theta'}\{L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I | T_{\mathbf{x}}(1) = k\}. \end{aligned} \quad (3)$$

But by the definition of $p_{\alpha,k}$, we have

$$P_{\Theta'}\{L_{B(\mathbf{x})}(\mathbf{x}) > -\alpha I | T_{\mathbf{x}}(1) = k\} \geq 1 - p_{\alpha,k}. \quad (4)$$

Summarizing (2)–(4), the proof of the lemma is complete. \square

IV. PROOFS OF THEOREMS 1 AND 2

Proof of Theorem 1: There are two cases. If $P_{\Theta'}\{T_{\mathbf{x}}(1) \leq c_n\} < a_n$, then by Markov's inequality, we have

$$E_{\Theta'}[T_{\mathbf{x}}(1)] \geq c_n(1 - a_n).$$

If, on the other hand, $P_{\Theta'}\{T_{\mathbf{x}}(1) \leq c_n\} \geq a_n$, then

$$\begin{aligned} P_{\Theta}\{T_{\mathbf{x}}(1) \leq c_n\} &= \sum_{k=1}^{c_n} P_{\Theta}\{T_{\mathbf{x}}(1) = k\} \\ &\geq \sum_{k=1}^{c_n} P_{\Theta'}\{T_{\mathbf{x}}(1) = k\} e^{-\alpha I} (1 - p_{\alpha,k}) \\ &\quad \text{(by Lemma 1, whenever } \alpha > c_n) \\ &\geq \Lambda_{\alpha, c_n} P_{\Theta'}\{T_{\mathbf{x}}(1) \leq c_n\} \\ &\geq a_n \Lambda_{\alpha, c_n}. \end{aligned}$$

Thus

$$P_{\Theta}\{T_{\mathbf{x}}(2) > n - c_n\} = P_{\Theta}\{T_{\mathbf{x}}(1) \leq c_n\} \geq a_n \Lambda_{\alpha, c_n}$$

and by Markov's inequality

$$E_{\Theta}[T_{\mathbf{x}}(2)] \geq (n - c_n) a_n \Lambda_{\alpha, c_n}.$$

Therefore

$$\begin{aligned} & \max(R_n(\Theta), R_n(\Theta')) \\ &= \Delta \max(E_{\Theta}[T_{\mathbf{x}}(2)], E_{\Theta'}[T_{\mathbf{x}}(1)]) \\ &\geq \Delta \min(c_n(1 - a_n), (n - c_n)a_n \Lambda_{\alpha, c_n}) \end{aligned}$$

and the theorem is proved. \square

Proof of Theorem 2: After time n , the regret under configuration Θ is

$$R_n(\Theta) = \Delta E_{\Theta}[T_{\mathbf{x}}(2)].$$

For any $\alpha > n$, we have

$$\begin{aligned} & \max(R_n(\Theta), R_n(\Theta')) \\ &\geq \frac{R_n(\Theta) + R_n(\Theta')}{2} \\ &= \Delta \frac{E_{\Theta}[T_{\mathbf{x}}(2)] + E_{\Theta'}[T_{\mathbf{x}}(1)]}{2} \\ &= \frac{\Delta}{2} \sum_{i=1}^n (P_{\Theta}\{T_{\mathbf{x}}(2) \geq i\} + P_{\Theta'}\{T_{\mathbf{x}}(1) \geq i\}) \\ &\geq \frac{\Delta}{2} \sum_{i=1}^n (P_{\Theta}\{T_{\mathbf{x}}(2) \geq i\} \\ &\quad + P_{\Theta}\{T_{\mathbf{x}}(1) \geq i\})e^{-\alpha i} (1 - p_{\alpha, k}) \quad (\text{by Lemma 1}) \\ &\geq \frac{\Delta}{2} \Lambda_{\alpha, n} \sum_{i=1}^n (P_{\Theta}\{T_{\mathbf{x}}(2) \geq i\} + P_{\Theta}\{T_{\mathbf{x}}(1) \geq i\}) \\ &= \frac{\Delta}{2} \Lambda_{\alpha, n} (E_{\Theta}[T_{\mathbf{x}}(2)] + E_{\Theta}[T_{\mathbf{x}}(1)]) \\ &= \frac{n\Delta}{2} \Lambda_{\alpha, n} \end{aligned}$$

and the proof is complete.

REFERENCES

- [1] R. Agrawal, M. Hedge, and D. Teneketzis, "Asymptotically efficient adaptive allocation rules for the multiarmed bandit problem with switching cost," *IEEE Trans. Automat. Contr.*, vol. 33, pp. 899–906, Oct. 1988.
- [2] R. Agrawal, D. Teneketzis, and V. Anantharam, "Asymptotically efficient adaptive allocation schemes for controlled i.i.d. processes: Finite parameter space," *IEEE Trans. Automat. Contr.*, vol. 34, pp. 258–266, Mar. 1989.
- [3] V. Anantharam, P. Varaiya, and J. Walrand, "Asymptotically efficient allocation rules for the multiarmed bandit problem with multiple plays—Part I: i.i.d. rewards," *IEEE Trans. Automat. Contr.*, vol. 32, pp. 968–976, Nov. 1987.
- [4] —, "Asymptotically efficient allocation rules for the multiarmed bandit problem with multiple plays—Part II: Markovian rewards," *IEEE Trans. Automat. Contr.*, vol. 32, no. 11, pp. 977–982, Nov. 1987.
- [5] D. A. Berry and B. Fristedt, *Bandit Problems*. New York: Chapman and Hall, 1985.
- [6] J. C. Gittins, *Multi-Armed Bandit Allocation Indices*. New York: Wiley, 1989.
- [7] T. L. Lai and H. Robbins, "Asymptotically efficient adaptive allocation rules," *Advances Appl. Math.*, vol. 6, pp. 4–22, 1985.
- [8] H. Robbins, "Some aspects of the sequential design of experiments," *Bull. Amer. Math. Soc.*, vol. 58, pp. 527–535, 1952.
- [9] W. R. Thompson, "On the likelihood that one unknown probability exceeds another in view of the evidence of two samples," *Biometrika*, vol. 25, pp. 275–294, 1933.

A Multivariable Bilinear Adaptive Controller with Decoupling Design

Xi Sun and Ming Rao

Abstract—This paper presents a new adaptive decoupling controller for multivariable bilinear systems with nondiagonal coefficient matrix of stochastic noises. The controller combines the feedforward control strategy with the generalized minimum variance approach and performs the decoupling of the control loop dynamically as well as in the steady state. It can control the bilinear systems whose linear part is not necessarily of minimum phase. The proof of global convergence for the algorithm is also provided. Simulation results demonstrate the effectiveness of the controller.

Index Terms—Adaptive control, bilinear systems, convergence, multi-variable decoupling.

I. INTRODUCTION

In many multivariable processes, there exists strong interaction between control loops. In such a case, it is important to consider a decoupling control strategy so as to improve the performance of the closed-loop systems. When the model parameters are unknown, a feasible approach is to adopt an adaptive decoupling scheme. Singh and Narendra [1] discussed adaptive decoupling and prior knowledge. Adaptive decoupling controllers which combine classical feedforward decoupling control with self-tuning control were proposed by McDermott and Mellichamp [2], Mahieddine and Morris [3], as well as Chai [4], [5]. In their papers, the cross-coupling terms of the systems were considered as measurable disturbances and the effects of their interactions were eliminated by the feedforward control. Those control schemes could guarantee static decoupling and approximate dynamic decoupling. Exact adaptive decoupling control was first developed using a decoupling precompensator by Wittenmark *et al.* [6]. The precompensator was designed in such a way that it is able to separate the controller design into several single-input/single-output (SISO) design problems. Nevertheless, all the above adaptive decoupling schemes are focused only on linear systems.

Bilinear systems are important because they are frequently encountered in chemical processes, biological systems, etc. [7]. Some adaptive control algorithms were developed for SISO bilinear systems [8]–[12]. For deterministic multivariable systems, Sen [13] described a model reference adaptive control law, and the asymptotic stability of the closed-loop system was obtained. But the control law requires that the delay time on the diagonal is shorter than the off-diagonal delay time.

This paper reconstructs the adaptive control algorithm for single variable bilinear system [11] into the one for multivariable bilinear systems. The decoupling design for multivariable bilinear systems is presented. It is shown that the closed-loop system is globally stable and asymptotically optimal in some sense, even for the bilinear systems whose open-loop nominal linear parts have stable or unstable zeros. The effectiveness of the controller is demonstrated through simulation results.

Manuscript received April 3, 1998; revised March 17, 1999. Recommended by Associate Editor, J. Farrell. This work was supported in part by Natural Sciences and Engineering Research Council of Canada and Natural Science Foundation of China.

The authors are with the Department of Chemical and Materials Engineering, The University of Alberta, Edmonton, AB T6G 2G6 Canada (e-mail: ming.rao@ualberta.ca).

Publisher Item Identifier S 0018-9286(00)04091-5.