Optimal Foraging and Learning

JOHN M. MCNAMARA

School of Mathematics, University of Bristol

AND

Alasdair I. Houston

Department of Zoology, University of Oxford

(Received 9 May 1985)

Optimal foraging theory usually assumes that certain key environmental parameters are known to a foraging animal, and predicts the animal's behaviour under this assumption. However, an animal entering a new environment has incomplete knowledge of these parameters. If the predictions of optimal foraging theory are to hold the animal must use a behavioural rule which both learns the parameters and optimally exploits what it has learnt. In most circumstances it is not obvious that there exists any simple rule which has both these properties. We consider an environment composed of well-defined patches of food, with each patch giving a smooth decelerating flow of food (Charnov, 1976). We present a simple rule which (asymptotically) learns about and optimally exploits this environment. We also show the rule can be modified to cope with a changing environment. We discuss what is meant by optimal behaviour in an unknown and possibly changing environment, using the simple rule we have presented for illustrative purposes.

Introduction

Optimal foraging theory is based on the assumption that natural selection will favour the foraging strategy that maximizes the forager's fitness (see Krebs *et al.*, 1983; Pyke *et al.*, 1977; Pyke, 1984 for reviews). What Krebs *et al.* (1983) call classical optimal foraging theory assumes that fitness is maximized by the maximization of the net rate of energetic gain. There are theoretical reasons why this net rate is not necessarily the appropriate currency (Caraco, 1980; McNamara & Houston, 1982), and there is evidence that considerations other than net rate determine foraging decisions (e.g. Caraco *et al.*, 1980; Caraco, 1983; Caraco & Lima, 1985). Nevertheless, we will work within the classical optimal foraging framework. The advantage of doing so is that it is easier to relate our arguments to previous work in this area. It is also worth pointing out that (a) net rate may sometimes be an appropriate currency and (b) the issues that we raise remain relevant whatever the currency for foraging may be.

The first models of optimal foraging assumed that the forager "knew" the parameters of the environment (e.g. Charnov, 1973; Pulliam, 1974). Oaten (1977) pointed out that rewards could provide a forager with information about its environment, but the forager was assumed to have knowledge of the possible environmental states. The forager is using information but not learning anything new about the environment as a whole. It can be argued that Krebs et al. (1978) were the first people to raise the issue of optimal learning. Great tits (Parus major) were given a series of choices between two feeding sites. Each site had a constant probability of giving the bird a food item if the bird chose it. At the start of a test, the reward probabilities were not known to the bird. This procedure corresponds to the two-armed bandit problem of decision theory. Problems such as this are difficult because they involve a combination of parameter estimation and the maximization of payoff. The optimal policy can be found by Bayesian decision theory in which the animal can be thought of as forming estimates of the reward probabilities (McNamara & Houston, 1980). There is, however, no need to assume that animals actually use such procedures. Simple rules can perform almost as well as a Bayesian decision maker (Houston et al., 1982).

Houston *et al.* (1982) consider simple learning rules within a framework based on optimization. In contrast to this approach, Ollason (1980) sees his learning rule as an alternative to an optimality analysis. Furthermore, he argues (p. 51) that if an animal is learning then it cannot be foraging optimally and if it is foraging optimally it cannot be learning. Ollason applied his learning rule to an environment of depleting patches. His simulations suggest that when there are several patch types his learning rule approaches the optimal behaviour. We use the same paradigm to explore the following issues: (a) is it always possible to learn the classical optimal strategy, i.e. the strategy which maximizes long term reward rate? (b) What does it mean to forage optimally in an unknown, and possible fluctuating, environment?

The Model Environment

We consider the patchy environment first analyzed by Charnov (1973, 1976). We assume that there are k patch types labelled E_1, E_2, \ldots, E_k . The proportion of patches of type E_i is α_i , so that the probability that the next patch visited is type E_i is always α_i . Patches are not revisited. On a patch of type E_i rewards are gained as a smooth flow: the rate after time t has

elapsed on a patch being $r_i(t)$. We assume that, for each *i*, $r_i(t)$ is a strictly decreasing function which tends to zero as *t* tends to ∞ . The travel time between patches, τ , is a random variable with a finite expectation $0 < E(\tau) < \infty$.

We assume that, on encountering this environment, the animal does not know k, α_i , r_i or τ .

The marginal value theorem (Charnov, 1976) says that the long term rate is maximized if the forager leaves each patch when the rate $r_i(t)$ drops to the maximum possible long-term rate. This rate will be denoted by γ^* . The theorem provides a rule that specifies the optimal behaviour if γ^* is known. We are assuming, however, that the environmental parameters, and hence γ^* , are not known. The forager is now faced with the problem that optimal behaviour (in its simplest sense) requires a knowledge of γ^* , but γ^* can only be achieved by behaving optimally. Despite this apparent circularity, we show in this section that it is always possible to learn γ^* in the environment that we are considering.

A general way in which behaviour can be represented is to view the state of the animal (or model animal) as determining behaviour. The consequences of behaviour in turn modify the state and hence determine future behaviour (e.g. Houston *et al.*, 1977). Within this framework, a rule can be said to learn about an environmental parameter when the value of a state variable converges to the parameter value in question. In the case of patch-use, the parameter that we will consider is the maximum possible rate γ^* . We now describe a simple rule that learns the value of γ^* .

The Learning Rule

The animal encounters patches sequentially. Let

 t_n = time spent on the *n*th patch

 τ_n = time spent travelling between the *n*th and *n* + 1th patch.

 $T_n = t_n + \tau_n.$

Thus $T_1 + \ldots + T_n$ is the time between arrival at the first patch and arrival at the n + 1st. Let

 R_n = reward obtained on the *n*th patch.

We consider a rule in which γ^* is recursively estimated. Initially constants R_0 and T_0 are chosen $(R_0, T_0 > 0)$. The ratio $\gamma_0 = R_0/T_0$ is the initial estimate for γ^* . The animal leaves the first patch encountered when the reward rate on this patch falls to γ_0 . On arrival at the next patch it forms a new estimate

 γ_1 for γ^* , and uses this to determine when to leave the second patch, and so on. We consider a rule which has the following detailed form.

(a) On arrival at the n + 1th patch (n = 0, 1, ...) the animal forms the estimate γ_n for γ^* where

$$\gamma_n = \frac{R_0 + R_1 + \ldots + R_n}{T_0 + T_1 + \ldots + T_n}.$$
 (1)

(b) The animal leaves the n+1th patch when the reward rate on this patch falls to γ_n .

To illustrate the properties of this learning rule, we first consider a simple environment in which all patches are the same. In this case it is possible to give a reasonably easy proof of results which are true in the more complex environments that we subsequently consider.

Learning in the Simplest Case

Consider an environment in which there is only one patch type with reward rate r(t), and the travel time τ is a constant.

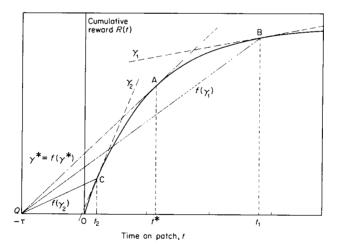


FIG. 1. A graphical illustration of the relationship between an estimate γ for γ^* and the long term reward rate $f(\gamma)$ which results from using this estimate on each patch encountered. The solid curve is the cumulative reward $R(t) = \int_0^t r(v) dv$ as a function of time t on patch. For any estimate γ , the rule is to leave each patch when r(t) falls to γ . The resulting rate $f(\gamma)$ is the cumulative reward R(t) divided by the total time $(t + \tau)$. Three values of γ are illustrated $(\gamma_1, \gamma^* \text{ and } \gamma_2)$. In each case the rate is the slope of the solid straight line that joins the curve to the point Q. When $\gamma = \gamma^*$, then $f(\gamma) = \gamma^*$. In the other two cases γ is indicated by the slope of the broken line.

Suppose that the animal uses the same estimate $\gamma > 0$ for γ^* on every patch it encounters. Let $f(\gamma)$ be the long term reward rate achieved by using this estimate. Figure 1 illustrates the dependence of $f(\gamma)$ on γ . It can be seen that

- (a) If $\gamma = \gamma^*$ then the animal leaves at time t^* and achieves rate $f(\gamma^*) = \gamma^*$.
- (b) If $\gamma = \gamma_1$ where $0 < \gamma_1 < \gamma^*$, the animal leaves at time t_1 and achieves rate $f(\gamma_1)$ where $\gamma_1 < f(\gamma_1) < \gamma^*$.
- (c) If $\gamma = \gamma_2$ where $\gamma_2 > \gamma^*$ then the animal leaves at time t_2 and achieves rate $f(\gamma_2)$ where $f(\gamma_2) < \gamma^*$.

Note that if $\gamma > r(0)$ then each patch is left immediately and hence $f(\gamma) = 0$. Figure 2 shows the full function f when $r(t) = e^{-t}$ and $\tau = 1$. In this case $\gamma^* = 0.31784$ and the optimal time on patch is $t^* = 1.1462$. As can be seen from the figure, conditions (a), (b) and (c) are all satisfied in this case.

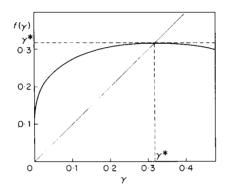


FIG. 2. The relationship between an estimate γ of γ^* and the long term reward rate $f(\gamma)$ which results from using this estimate on each patch encountered. The curve gives $f(\gamma)$ when $r(t) = e^{-t}$ and $\tau = 1.0$. The points of intersection of this curve with the straight line give solutions to the equation $f(\gamma) = \gamma$. Note that $\gamma = \gamma^*$ is the only positive solution to this equation.

The graphical results, illustrated in Fig. 1, are true in general and are established by analytical argument in Appendix 1. We summarise the main conclusions.

$$0 \le f(\gamma) \le \gamma^*$$
 for all $\gamma > 0$ (2)

$$f(\gamma^*) = \gamma^* \tag{3}$$

$$0 < \gamma < \gamma^* \Longrightarrow \gamma < f(\gamma). \tag{4}$$

In particular γ^* is the only positive root of the equation $f(\gamma) = \gamma$: i.e.

$$f(\gamma) = \gamma \quad \text{for } \gamma > 0 \Longrightarrow \gamma = \gamma^*.$$
 (5)

We now apply these results to our learning rule. Equations (2), (3) and (4) show that the estimates γ_n and γ_{n+1} satisfy the following relationships

$$\gamma_n > \gamma^* \Longrightarrow 0 < \gamma_{n+1} < \gamma_n. \tag{6}$$

$$\gamma_n = \gamma^* \Longrightarrow \gamma_{n+1} = \gamma^*. \tag{7}$$

$$0 < \gamma_n < \gamma^* \Longrightarrow \gamma_n < \gamma_{n+1} < \gamma^*. \tag{8}$$

These relationships are proved in Appendix 2. They imply the following result.

$$\gamma_n \to \gamma^* \text{ as } n \to \infty.$$
 (9)

The proof can be divided into two parts.

Part 1. Proof that γ_n tends to a limit Suppose that $0 < \gamma_0 < \gamma^*$. Then by equation (8)

$$\gamma_0 < \gamma_1 < \gamma_2 < \ldots < \gamma^*.$$

Thus γ_n tends to a limit as *n* tends to infinity.

Suppose $\gamma_0 = \gamma^*$. Then, by equation (7), $\gamma_n = \gamma^*$ for all *n*.

Finally, suppose that $\gamma_0 > \gamma^*$. Then, by equation (6), one of the following two possibilities must occur

(a) $\gamma_0 \ge \gamma_1 \ge \gamma_2 \ge \ldots \ge \gamma^*$, or

(b) $0 < \gamma_N < \gamma^*$ for some N.

If (a) occurs then γ_n must tend to a limit as *n* tends to infinity. If (b) occurs then by equation (8)

$$\gamma_N < \gamma_{N+1} < \gamma_{N+2} < \ldots < \gamma^*,$$

and again we have convergence.

In conclusion γ_n tends to a limit as *n* tends to infinity. We denote this limit by $\tilde{\gamma}$. It can be seen from the above proof that $\tilde{\gamma} > 0$.

Part 2. Proof that $\tilde{\gamma} = \gamma^*$

Since γ_n tends to a limit so do the rewards R_n and times T_n . We denote these limits by R and T respectively. By equation (1) γ_n must tend to R/T, so that $R = \tilde{\gamma}T$.

It can be shown that $f(\gamma)$ is a continuous function of γ for $\gamma > 0$. Thus $f(\gamma_n)$ tends to $f(\tilde{\gamma})$ as *n* tends to infinity. But $R_{n+1} = f(\gamma_n) T_{n+1}$, since $f(\gamma_n)$ is the average rate the animal would obtain by using the estimate γ_n on each patch. Thus taking the limit as *n* tends to infinity we find that $R = f(\tilde{\gamma}) T$. Comparing this with the equation $R = \tilde{\gamma}T$ we conclude that $f(\tilde{\gamma}) = \tilde{\gamma}$. Thus from condition (5) it can be seen that $\tilde{\gamma} = \gamma^*$.

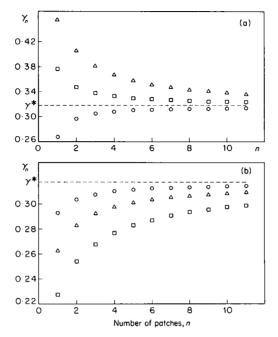


FIG. 3. The convergence of γ_n to γ^* when $r(t) = e^{-t}$ and $\tau = 1.0$. In this case $\gamma^* = 0.31784$. In (a) $\gamma_0 = 0.6 > \gamma^*$. The triangles correspond to $T_0 = 2$, $R_0 = 1.2$, the squares to $T_0 = 0.75$, $R_0 = 0.45$ and the circles to $T_0 = 0.01$, $R_0 = 0.006$. In (b) $\gamma_0 = 0.15 < \gamma^*$. The squares correspond to $T_0 = 2.5$, $R_0 = 0.375$, the triangles to $T_0 = 0.8$, $R_0 = 0.12$ and the circles to $T_0 = 0.01$, $R_0 = 0.0015$.

The convergence of γ_n to γ^* is illustrated in Fig. 3, from which the effect of variations in the choice of R_0 and T_0 can be seen. In Fig. 3(a) γ_0 is greater than γ^* and in Fig. 3(b) γ_0 is less than γ^* . When $R_0 = 0.006$, $T_0 = 0.01$, it can be seen that $\gamma_0 > \gamma^*$ and $\gamma_1 < \gamma^*$. In all the other examples illustrated, convergence is monotone.

Learning in the General Case

When there is more than one patch type present (or when τ varies) T_n and R_n are random variables and the sequence $\gamma_0, \gamma_1, \gamma_2, \ldots$ is a stochastic process. We illustrate this with three stochastic environments which are all simple modifications of the environment in which each patch has reward rate $r(t) = e^{-t}$.

In each environment there are two patch types E_1 and E_2 , with parameters k_1 and k_2 respectively. On a patch with parameter k_i (i = 1, 2) we have

 $r_i(t) = k_i \exp(-t)$. The two types are present in equal proportions so that $\alpha_1 = \alpha_2 = 0.5$. The travel time τ is again equal to 1. The three environments are described in Table 1 together with the optimal stay times on each patch type. The table also gives the resulting long term average rate on each patch type (rewards on patch divided by time on patch plus τ). For each environment parameters have been chosen so that the maximum long term average rate, γ^* , is again equal to 0.31784.

| Environment | Patch type | k | <i>t</i> * | Rate |
|-------------|-----------------------------------|--------|------------|--------|
| A | E_1 | 0.8 | 0.9231 | 0.2507 |
| | | 1.1823 | 1.3137 | 0.3736 |
| В | $\tilde{E_1}$ | 0.6 | 0.6354 | 0.1725 |
| | $\dot{E_{2}}$ | 1.3277 | 1.4297 | 0.4156 |
| С | $E_2 \\ E_1 \\ E_2 \\ E_1 \\ E_2$ | 0.4 | 0.2299 | 0.0668 |
| | Ė, | 1.4203 | 1.4971 | 0.4415 |

TABLE 1

To illustrate the stochastic nature of the process, consider Environment C and suppose that the initial constants are

$$T_0 = 1$$
 and $R_0 = 0.31784 \equiv \gamma^*$,

so that $\gamma_0 = \gamma^*$. The possible values of γ_1 , γ_2 and γ_3 are given in Fig. 4. As can be seen there are two possible values of γ_1 , four possible values of γ_2 , and eight possible values of γ_3 . Each time that a good patch (type E_2) is encountered, γ_n increases; each time that a poor patch is encountered, γ_n decreases. Because $\alpha_1 = \alpha_2 = 0.5$, the eight possible values of γ_3 are equally likely.

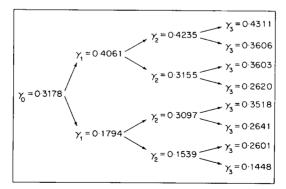


FIG. 4. Possible values of γ_1 , γ_2 and γ_3 for Environment C when $T_0 = 1$ and $R_0 = 0.31784 = \gamma^*$.

It is important to note that the order in which patch types are visited influences γ_n . For example the sequence E_2, E_2, E_1 results in $\gamma_3 = 0.3606$, whereas the sequence E_1 , E_2 , E_2 results in $\gamma_3 = 0.3518$. The reason for this effect is that the patch types that have been visited determine the current estimate of γ , which in turn determines the time spent on the current patch. Thus the time spent on a patch type and the rewards obtained from it are not fixed, but depend on previous experience. As a result, even when $\gamma_0 = \gamma^*$ it is not clear that γ_n tends to γ^* as *n* tends to infinity. The difficulty is that even though the proportion of patches of type *i* that are visited tends to α_i . as the number of patches visited tends to infinity, the order in which patches are visited is also relevant. It is therefore not possible to extend the proof given above by simply averaging. By approaching the problem in a different way, it is possible however to prove that equation (9) holds with probability one. This result follows from an application of a general result on Markov renewal processes given in McNamara (1985). Strictly speaking the general theory presented there only applies to the case where each patch visited eventually runs out of food, so that $r_i(t) = 0$ for t sufficiently large. The theory can however be extended in various ways to cover other cases.

Ollason (1980) considers a rule which has some similarities to the rule we consider. One major difference is that the past is discounted by his rule. The degree of discounting is measured by a parameter k which tends to infinity as memory stretches further back into the past. He considers the special case of rewards of the form $r(t) = ae^{-bt}$ in detail. In this case his rule has the property that, when all patch are identical, the time spent on patch tends to a limit, t_c , as time in the environment tends to infinity. As k tends to infinity t_c tends to the optimal time on patch. When patches differ his simulations suggest that behaviour tends towards optimal as time in the environment and k tend to infinity.

This suggests that setting $k = \infty$ might give a rule which, like our rule, is asymptotically optimal (see below). As can be seen from his equation (1) this is not the case. Setting $k = \infty$ here leads to a rule under which a patch is never left until r(t) = 0. Such a rule would perform very badly.

Rates of Convergence

Clearly it is advantageous to obtain accurate estimates for γ^* as quickly as possible. To measure how quickly the estimate γ_n tends to γ^* we introduce the error σ_n^2 defined by

$$\sigma_n^2 = E\{(\gamma_n - \gamma^*)^2\};$$

i.e. σ_n^2 is the mean square error in the estimate γ_n of γ^* . We use σ_n^2 to

measure estimation error because the loss in reward rate that results from using the wrong γ to determine when to leave a patch is approximately proportional to $(\gamma - \gamma^*)^2$ for small $|\gamma - \gamma^*|$. A justification of this claim is given in Appendix 3.

If $\gamma_0 = \gamma^*$ then $\sigma_n^2 = 0$ for all *n* in a non-stochastic environment. Thus in order to investigate the effect of stochasticity on σ_n^2 , we set $T_0 = 1$ and $R_0 = \gamma^*$, so that $\gamma_0 = \gamma^*$. The resulting mean square error in the three environments A, B and C is illustrated in Fig. 5, which shows that the mean square error increases with increasing variability in the environment.

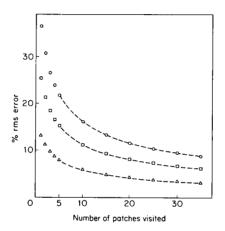


FIG. 5. The effect of environmental variability on the error with which γ^* is estimated. The root mean square error σ_n is plotted as a percentage of γ^* for (from top to bottom) the three environments C, B and A as a function of the number, *n*, of patches visited. The environments are described in Table 1, which shows that C is the most variable and A the least variable. The figure shows that the error increases with increasing environmental variability.

To illustrate the additional effect of a variable travel time, some values of the percentage root mean square error under variable τ are given in Table 2. The values come from Environment B, modified so that τ has mean 1 but a uniform distribution over an interval centred on 1. The table shows that the mean square error increases with increasing variability in τ .

| TABLE 2 | 2 |
|---------|---|
|---------|---|

Learning γ^* when the travel time is variable (Environment B)

| % rms error | <i>n</i> = 2 | <i>n</i> = 4 | n = 8 | n = 20 | n = 50 |
|---|--------------|--------------|-------|--------|--------|
| $\tau = 1$ | 21.2 | 16.5 | 12.3 | 8.1 | 5.2 |
| $\tau \sim U(\frac{1}{2}, \frac{3}{2})$ | 22.7 | 17.8 | 13.2 | 8.6 | 5.5 |
| $\tau \sim U(0,2)$ | 27.6 | 21.5 | 15.8 | 10.2 | 6.5 |

These results illustrate a phenomenon which we might expect to hold generally for any rule which estimates γ^* : the more variability in the environment, the slower an estimate tends to converge to γ^* .

It is clear from Fig. 3 that the rate at which γ_n tends to γ^* also depends on R_0 and T_0 . Mean square error is reduced if the initial estimate γ_0 is close to γ^* . However, it is not always possible to ensure this since γ^* is a priori unknown. For given $\gamma_0 = R_0/T_0$ Fig. 3 illustrates that σ_n^2 also depends on the magnitudes of R_0 and T_0 . When all patches are the same small R_0 and T_0 ensure rapid convergence. (Although, as Fig. 3 illustrates, if $\gamma_0 > \gamma^*$ and R_0 and T_0 are very small γ_1 may overshoot γ^* .) However, if patches differ, this is no longer true. As Fig. 4 illustrates random fluctuations can change the estimate of γ^* dramatically. The smaller the values of R_0 and T_0 the more sensitive a rule is to these fluctuations. As a consequence if one is fairly sure that the initial estimate γ_0 is close to γ^* it is best to choose R_0 and T_0 large so that estimates will be insensitive to short term fluctuations.

It can thus be seen that the initial estimates R_0 and T_0 together define a sort of prior mean and variance for the unknown parameter γ^* . The ratio $\gamma_0 = R_0/T_0$ acts as the prior mean, while the magnitude of the parameters R_0 and T_0 play the role of prior variance, with variance decreasing with increasing R_0 and T_0 .

Optimal Learning

The discussion of convergence has led us to the central issue of optimality in the context of learning. As we have already mentioned, Ollason (1980) argues that learning is incompatible with optimal foraging, and vice versa. While it is clearly true that an animal (or rule) that has to learn about the environment cannot always maximize reward rate, it does not seem reasonable to conclude that a rule for learning cannot be optimal in some broader sense. One such sense is that of asymptotic optimality, which has already been described in the context of learning by Houston *et al.* (1982). A rule is asymptotically optimal if its performance tends to the optimal performance as experience tends to infinity. In other words, the rule always learns the optimal policy in the long run. This criterion is more or less the one used by Harley (1981). The results described above show that there exists an asymptotically optimal rule for a broad class of patchy environments.

It can be objected, however, that asymptotic optimality is not necessarily a very good criterion, in that it assumes that performance will be assessed over an unlimited period. Rules that are asymptotically optimal may perform badly over limited periods (Houston *et al.*, 1982). The appropriate currency depends on the details of an animal's environment. If the animal has a fixed period in a given environment, then the total expected reward may be the best criterion. If the animal has a constant probability of being forced to stop foraging, then some form of discounted rewards should be the criterion (see Houston *et al.*, 1982).

Establishing the optimality criterion does not, however, suffice to enable the optimal learning rule to be specified. This can be seen from our discussion of how our rule depends on R_0 and T_0 , where it was seen that the best choice of these parameters depends on the type of environment likely to be encountered. In general one needs to specify the range of environments which are possible, and to specify the probability of encountering an environment of any given type. In other words one needs to specify some prior distribution on the set of possible environments.

Finally, one needs to specify what information a rule is allowed to use.

Coping with a Changing Environment

So far, we have confined our attention to a constant, but possibly stochastic, environment. The principle advantage of learning is, however, the fact that it enables an animal to respond to changes in the environment. Many people have constructed learning rules that are based on some exponentially weighted average of past experience (Harley, 1981; Killeen, 1982; Lester, 1984; Ollason, 1980; see Kacelnik *et al.*, in press, for a review). The basic idea of exponential weighting is that more importance is given to the recent as opposed to the remote past. We now consider a modified rule that involves exponential weighting, and investigate how well this rule performs in both changing and non-changing environments. This illustrates the conflicting pressures on a learning rule and leads us to a consideration of what it means to be optimal in a changing invironment.

Following the notation used above, R_n is the reward obtained on the *n*th patch and T_n is the time on the *n*th patch plus the travel time to the next patch. Like our previous rule, the rule forms an estimate of γ^* using previous rewards and times and leaves a patch when the rate on the patch falls to the current estimate of γ . The previous rule formed the following estimate of γ_n

$$\gamma_n = \frac{R_n + \ldots + R_0}{T_n + \ldots + T_0}.$$

In this estimate, all previous rewards and times have equal weight. The rule that we now consider is based on the following estimate of γ_n

$$\gamma_n = \frac{R_n + e^{-\alpha T_n} R_{n-1} + e^{-\alpha (T_n + T_{n-1})} R_{n-2} + \ldots + e^{-\alpha (T_n + \ldots + T_1)} R_0}{T_n + e^{-\alpha T_n} T_{n-1} + e^{-\alpha (T_n + \ldots + T_{n-1})} T_{n-2} + \ldots + e^{-\alpha (T_n + \ldots + T_1)} T_0}$$

In this equation, α is the exponential weighting factor ($\alpha > 0$) (the previous rule corresponds to $\alpha = 0$). A simple way to express this rule is to call the numerator Y_n and the denominator X_n . It can then be seen that $\gamma_n = Y_n/X_n$ where

$$Y_n = R_n + \mathrm{e}^{-\alpha T_n} Y_{n-1}$$

and

$$X_n = T_n + \mathrm{e}^{-\alpha T_n} X_{n-1}.$$

To investigate the properties of this rule we first consider an unchanging environment in which all patches are the same and the travel time is constant. In such an environment it can again be shown that γ_n tends to γ^* as *n* tends to infinity. The proof is a simple modification of the proof given for the original rule. As before one can establish equations (6), (7) and (8) and use these equations to show that γ_n must converge to some limit, and one can again use equation (5) to show that this limit must be γ^* . Thus the modified rule is asymptotically optimal for any value of the exponential weighting α . The rule proposed by Ollason (1980) does not have this property for any value of his weighting factor k.

Under this rule the quantities Y_n and X_n tend to equilibrium values as the number of patches visited tends to infinity. The equilibrium values depend on α and the environment. The rate at which equilibrium is approached depends on α and increases with increasing α .

The point of using a rule based on some weighting or discounting of past experience is that a rapid response to change is possible. To illustrate the effect of the weighting factor α , we have calculated the response of the rule to a change in patch quality. We assume that the rule has reached equilibrium in an environment in which the rate on all patches is $r(t) = k \exp(-t)$ and k = 1. The value of k then changes to a new value which is the same for all patches. Figure 6 shows how the estimate of γ^* depends on α and on the number of patches that have been visited. When α is large the previous environment is forgotten quickly and convergence to the new γ^* is rapid.

Figure 6 suggests that a large value of α is desirable, but this is not true if the environment is stochastic. When there is no exponential weighting, all patches are represented equally in the estimate of γ^* . The introduction of a weighting factor means that recent patches predominate. No matter how many patches have been visited the estimate γ_n is subject to the influence of the patches which happen to have been visited recently. Consequently, when the environment is stochastic γ_n will not tend to a limit, but will forever fluctuate. One can use the mean square error to measure the size of these fluctuations. When there is no weighting ($\alpha = 0$) the mean square

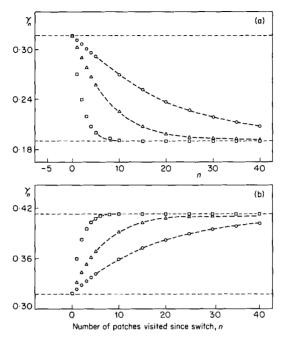


FIG. 6. Response of the rule to a change in the environment. Initially $r(t) = e^{-t}$ for all patches and the estimate of γ^* is its true value $\gamma^* = 0.31784$. After n = 0 the environment changes so that $r(t) = k e^{-t}$ for all patches. In (a) k = 0.6 so that the new γ^* is 0.1907. In (b) k = 1.3 so that $\gamma^* = 0.4132$. In each case squares refer to $\alpha = 0.25$ ($t_r = 4$), triangles to $\alpha = 0.06666$ ($t_r = 16$) and circles to $\alpha = 0.025$ ($t_r = 40$). $\tau = 1$ throughout.

error tends to zero, but when there is a non-zero weighting factor α the mean square error tends to a positive value. The dependence of this value on the response time $t_r = 1/\alpha$ for environments A, B, and C is shown in Fig. 7. The shorter the value of t_r , the more weight is put on the very recent past, and hence the more sensitive γ_n is to environmental stochasticity. Thus mean square error is large for small t_r and decreases as t_r increases.

It can be seen from this discussion that there are conflicting pressures on the choice of α . The need to respond to a change in the environment favours a small value of α , but such a value may result in inappropriate responses to runs of good or bad luck.

To determine the best value of α , we need to say what it means for a learning rule to be optimal in a changing environment. We have already said that in a stochastic environment it is necessary to specify a prior distribution on environments i.e. to specify the possible environments and their probability of occurrence. In a changing environment this information

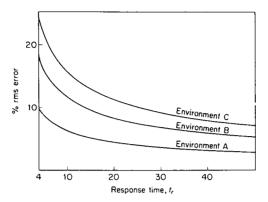


FIG. 7. The long term equilibrium value of the % rms error plotted as a function of $t_r = 1/\alpha$ for environments A, B and C.

is still required but in addition we must specify how frequently changes occur and the probability distribution of environments after a change. If the state of the environment is regarded as a vector-valued stochastic process, then the above requirements amount to having knowledge of the transition probabilities that govern the evolution of the stochastic process.

For the rule that we have been considering the best value of α depends on the probability that the environment will change. If a change is likely then a large α , and hence small response time, are likely to be advantageous. This advantage is offset, however, by the stochasticity of the environment. The interaction between these conflicting demands is quite complex.

Discussion

Optimal foraging theory usually assumes that certain environmental parameters are known. Ollason (1980) stressed that animals have to learn about their environment and discussed this problem in the context of patch use. The marginal value theorem (Charnov, 1976) specifies optimal behaviour in terms of the maximum possible rate γ^* , but as Ollason pointed out, if a predator does not know γ^* , then it can only experience γ^* by behaving optimally. Thus it is not clear that an animal can learn to behave optimally. We have presented a simple rule and shown that the rule learns γ^* under a wide range of conditions.

Ollason (1980) argues that because a learning animal does not spend all its time foraging at the maximum rate γ^* it is not foraging optimally. While this is obviously true in terms of classical foraging theory, it is inappropriate to use such an optimality criterion when the environmental parameters are unknown. We argue in this paper that the optimality criterion must take account of the fact that an animal has to learn. It is still possible to ask if an animal is foraging as well as possible, given that it starts without complete knowledge of its environment. We have described the ingredients that are needed to formulate an optimality criterion for behaviour in a constant but unknown environment and in a changing environment. Although this involves more ingredients than are required in classical foraging theory, it is still possible for optimal behaviour to be well-defined.

Our simulations suggest that it may be difficult to learn when the environment is highly variable. The basic model involves only variability in patch types; even then learning is slow. When variability of travel times is added, learning is slower. The foraging environment of most animals is presumably more variable than any of the cases that we have considered. Consequently, it may be very difficult for animals to learn enough about their environments to forage optimally in the terms of classical foraging models.

A.I.H. was supported by a Science and Engineering Research Council grant to J. M. McNamara and J. R. Krebs.

REFERENCES

CARACO, T. (1980). Ecology 61, 119.

- CARACO, T. (1983). Behav. Ecol. Sociobiol. 12, 63.
- CARACO, T. & LIMA, S. L. (1985). Anim. Behav. 33, 216.
- CARACO, T., MARTINDALE, S. & WHITTAM, T. S. (1980). Anim. Behav. 28, 820.
- CHARNOV, E. L. (1973). Optimal Foraging—some theoretical considerations. Ph.D. thesis, University of Washington.
- CHARNOV, E. L. (1976). Theor. Pop. Biol. 9, 129.
- HARLEY, C. B. (1981). J. theor. Biol. 89, 611.
- HOUSTON, A. I., HALLIDAY, T. R. & MCFARLAND, D. J. (1977). Med. Biol. Eng. Comput. 15, 49.
- HOUSTON, A. I., KACELNIK, A. & MCNAMARA, J. M. (1982). In: Functional Ontogeny (McFarland, D. J. ed.). London: Pitman.
- KACELNIK, A., KREBS, J. R. & ENS, B. (in press). In: *The Quantitative Analysis of Behaviour*, Vol. 6 (Commons, M. L., Shettleworth, S. J. & Kacelnik, A. eds). New York: Lawrence Erlbaum.
- KILLEEN, P. R. (1982). In: Nebraska Symposium on Motivation (Bernstein, D. J. ed.). Lincoln: University of Nebraska Press.
- KREBS, J. R., KACELNIK, A. & TAYLOR, P. (1978). Nature 275, 27.
- KREBS, J. R., STEPHENS, D. W. & SUTHERLAND, W. J. (1983). In: Perspectives in Ornithology (Brush, A. H. & Clark, G. A. eds). Cambridge: Cambridge University Press.
- LESTER, N. P. (1984). Behaviour 89, 175.
- MCNAMARA, J. M. (1985). J. Appl. Prob. 22, 324.
- MCNAMARA, J. M. & HOUSTON, A. I. (1980). J. theor. Biol. 85, 673.
- MCNAMARA, J. M. & HOUSTON, A. I. (1982). In: Functional Ontogeny (McFarland, D. J. ed.). London: Pitman.
- OATEN, A. (1977). Theor. Pop. Biol. 12, 263.
- OLLASON, J. G. (1980). Theor. Pop. Biol. 18, 44.

PULLIAM, H. R. (1974). Am. Nat. 108, 59. PYKE, G. H. (1984). Ann. Rev. Ecol. Syst. 15, 523. PYKE, G. H., PULLIAM, H. R. & CHARNOV, E. L. (1977). Quart. Rev. Biol. 52, 137.

APPENDIX 1

Properties of $f(\gamma)$

Let $\gamma > 0$. Since r(t) is monotonic decreasing and tends to zero as t tends to infinity there exists $t \equiv t(\gamma)$ such that

$$r(t) > \gamma \qquad 0 \le t < t(\gamma) \tag{A1.1}$$

$$r(t(\gamma)) \le \gamma. \tag{A1.2}$$

Note that, since r(t) is a decreasing function of t, $t(\gamma)$ is a decreasing function of γ .

We now define R by

$$R(t) = \int_0^t r(v) \, \mathrm{d}v \qquad t \ge 0, \tag{A1.3}$$

and define γ^* by

$$\gamma^* = \sup_{t \ge 0} \frac{R(t)}{(t+\tau)}.$$
 (A1.4)

Finally, we define

$$f(\gamma) = \frac{R(t(\gamma))}{t(\gamma) + \tau} \qquad \gamma > 0. \tag{A1.5}$$

We show that f satisfies equations (2), (3) and (4).

Equation (2) holds by equations (A1.4) and (A1.5), and equation (3) is just a restatement of the Marginal Value Theorem. Thus it only remains to establish equation (4). Let γ satisfy $0 < \gamma < \gamma^*$. Then, since $t(\gamma)$ is monotone decreasing, we have $t(\gamma) \ge t(\gamma^*)$. We thus have

$$R(t(\gamma)) = \int_0^{t(\gamma)} r(v) \, \mathrm{d}v = R(t(\gamma^*)) + \int_{t(\gamma^*)}^{t(\gamma)} r(v) \, \mathrm{d}v.$$
(A1.6)

Now by equation (A1.5)

$$R(t(\gamma^*)) = (t(\gamma^*) + \tau)f(\gamma^*)$$

Thus by equation (3) and the assumption $\gamma < \gamma^*$ we have

$$R(t(\gamma^*)) > \gamma(t(\gamma^*) + \tau). \tag{A1.7}$$

Also by equation (A1.1)

$$\int_{t(\gamma^*)}^{t(\gamma)} r(v) \, \mathrm{d}v \ge \int_{t(\gamma^*)}^{t(\gamma)} \gamma \, \mathrm{d}v = \gamma(t(\gamma) - t(\gamma^*)). \tag{A1.8}$$

Thus from equations (A1.6), (A1.7) and (A1.8) we have

$$R(t(\gamma)) > \gamma(t(\gamma) + \tau).$$

Equation (4) then follows from the definitions of $f(\gamma)$ (equation A1.5).

APPENDIX 2

Relationship of γ_{n+1} to γ_n

We prove equations (6), (7) and (8).

The learning rule prescribes that the n+1th patch should be left when the reward rate falls to γ_n . Thus in the terminology of Appendix 1

$$R_{n+1} = R(t(\gamma_n))$$

and

$$T_{n+1} = \tau + t(\gamma_n).$$

Thus by equation (A1.5)

$$R_{n+1} = f(\gamma_n) T_{n+1}.$$
 (A2.1)

For convenience we set

$$X_n = \sum_{i=0}^{n} T_i$$
 and $Y_n = \sum_{i=0}^{n} R_i$

so that

$$\gamma_n = \frac{Y_n}{X_n}.$$
 (A2.2)

We also have

$$\gamma_{n+1} = \frac{Y_n + R_{n+1}}{X_n + T_{n+1}},$$

and hence

$$\gamma_{n+1} = \frac{\gamma_n X_n + f(\gamma_n) T_{n+1}}{X_n + T_{n+1}}$$
(A2.3)

by equations (A2.1) and (A2.2).

Suppose that $\gamma_n > \gamma^*$. Then $f(\gamma_n) \le \gamma^*(<\gamma_n)$ by equation (2), and hence $\gamma_n X_n + f(\gamma_n) T_{n+1} < \gamma_n (X_n + T_{n+1})$. Equations (A2.3) then shows that $\gamma_{n+1} < \gamma_n$. We also have $\gamma_{n+1} > 0$ since X_n and γ_n are positive. This proves equation (6).

Suppose $\gamma_n = \gamma^*$. The $f(\gamma_n) = \gamma^*$ by equation (3) and hence $\gamma_{n+1} = \gamma^*$ by equation (A2.3). This proves equation (7).

Finally, suppose $0 < \gamma_n < \gamma^*$. Then $f(\gamma_n) > \gamma_n$ by equation (4), and $\gamma^* \ge f(\gamma_n)$ by equation (2). Thus

$$\gamma_n(X_n + T_{n+1}) < \gamma_n X_n + f(\gamma_n) T_{n+1} < \gamma^*(X_n + T_{n+1}).$$

Equation (8) then follows from this equation and equation (A2.3).

APPENDIX 3

Mean Square Error

The function $f(\gamma)$ has a maximum at $\gamma = \gamma^*$. Thus

$$f'(\gamma^*) = 0 \tag{A3.1}$$

and

$$k \equiv -\frac{1}{2} f''(\gamma^*) \ge 0.$$
 (A3.2)

Performing a Taylor series expansion of $f(\gamma)$ about $\gamma = \gamma^*$ we obtain

$$f(\gamma) = f(\gamma^*) + (\gamma - \gamma^*)f'(\gamma^*) + \frac{1}{2}(\gamma - \gamma^*)^2 f''(\gamma^*) + 0((\gamma - \gamma^*)^3).$$

Thus by equations (A3.1) and (A3.2)

$$f(\gamma^*) - f(\gamma) = k(\gamma - \gamma^*)^2 + O((\gamma - \gamma^*)^3).$$

This equation shows that the loss in long term reward rate which results from using the estimate γ rather than γ^* is proportional to $(\gamma - \gamma^*)^2$ for small $|\gamma - \gamma^*|$.