

pulling  $M_\lambda$ , so it gives up a reward of  $\lambda$  each time.) With just one arm, the only choice is to whether to pull again or to stop. If you pull the first arm  $T$  times (i.e., at times  $0, 1, \dots, T - 1$ ) we say that the **stopping time** is  $T$ .

Stopping time

Going back to our version with  $M$  and  $M_\lambda$ , let's assume that after  $T$  pulls of the first arm, an optimal strategy eventually pulls the second arm for the first time. Since no information is gained from this move (we already know the payoff will be  $\lambda$ ), at time  $T + 1$  we will be in the same situation and thus an optimal strategy must make the same choice.

Equivalently, we can say that an optimal strategy is to run arm  $M$  up to time  $T$  and then switch to  $M_\lambda$  for the rest of time. It's possible that  $T = 0$  if the strategy chooses  $M_\lambda$  immediately, or  $T = \infty$  if the strategy never chooses  $M_\lambda$ , or somewhere in between. Now let's consider the value of  $\lambda$  such that an optimal strategy is *exactly indifferent* between (a) running  $M$  up to the best possible stopping time and then switching to  $M_\lambda$  forever, and (b) choosing  $M_\lambda$  immediately. At the tipping point we have

$$\max_{T>0} E \left[ \left( \sum_{t=0}^{T-1} \gamma^t R_t \right) + \sum_{t=T}^{\infty} \gamma^t \lambda \right] = \sum_{t=0}^{\infty} \gamma^t \lambda,$$

which simplifies to

$$\lambda = \max_{T>0} \frac{E \left( \sum_{t=0}^{T-1} \gamma^t R_t \right)}{E \left( \sum_{t=0}^{T-1} \gamma^t \right)}. \tag{17.15}$$

This equation defines a kind of “value” for  $M$  in terms of its ability to deliver a stream of timely rewards; the numerator of the fraction represents a utility while the denominator can be thought of as a “discounted time,” so the value describes the maximum obtainable utility per unit of discounted time. (It's important to remember that  $T$  in the equation is a stopping time, which is governed by a rule for stopping rather than being a simple integer; it reduces to a simple integer only when  $M$  is a deterministic reward sequence.) The value defined in Equation (17.15) is called the **Gittins index** of  $M$ .

Gittins index

The remarkable thing about the Gittins index is that it provides a very simple optimal policy for any bandit problem: *pull the arm that has the highest Gittins index, then update the Gittins indices.* Furthermore, because the index of arm  $M_i$  depends only on the properties of that arm, an optimal decision on the first iteration can be calculated in  $O(n)$  time, where  $n$  is the number of arms. And because the Gittins indices of the arms that are not selected remain unchanged, each decision after the first one can be calculated in  $O(1)$  time.



### 17.3.1 Calculating the Gittins index

To get more of a feel for the index, let's calculate the value of the numerator, denominator, and ratio in Equation (17.15) for different possible stopping times on the deterministic reward sequence  $0, 2, 0, 7.2, 0, 0, 0, \dots$ :

$T$	1	2	3	4	5	6
$R_t$	0	2	0	7.2	0	0
$\sum \gamma^t R_t$	0.0	1.0	1.0	1.9	1.9	1.9
$\sum \gamma^t$	1.0	1.5	1.75	1.875	1.9375	1.9687
ratio	0.0	0.6667	0.5714	1.0133	0.9806	0.9651

Clearly, the ratio will decrease from here on, because the numerator remains constant while the denominator continues to increase. Thus, the Gittins index for this arm is 1.0133, the