**CHAPTER 10**

# OPTIMIZING A SCALAR FUNCTION

Optimizing scalar functions arises in a variety of settings, such as choosing the price of a product, the amount of memory in a computer, the temperature of a chemical process or the diameter of the tube in an aerosol gun. It can even be used for optimizing a tunable parameter in a learning policy! Important special cases are functions which are unimodular, as depicted in Figure 10.1, where there is a single local maximum, or concave. In this chapter, we review some specialized algorithms designed for this particular problem class.

We begin our presentation using an unusual problem setting for this volume, which is an unknown but deterministic function which can be measured perfectly. We then transition back to our more familiar setting of noisy functions.

## 10.1 DETERMINISTIC EXPERIMENTS

We are going to consider three settings where we wish to find the best estimate of the optimum of a deterministic, unimodular function within a given experimental budget. The settings include

- Differentiable functions.

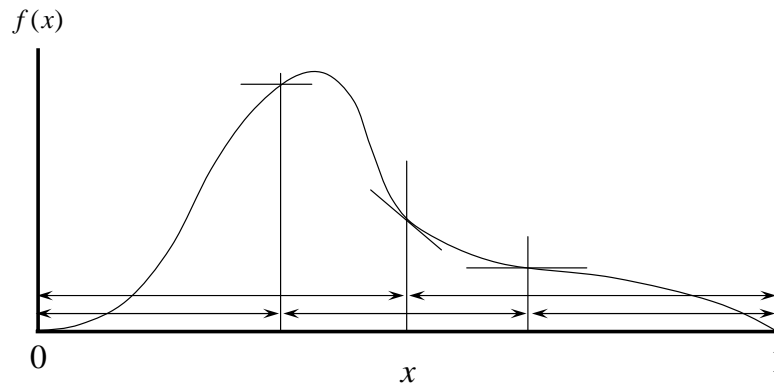- Nondifferentiable functions, with a finite experimental budget.

**Figure 10.1** A unimodular function, showing the slope at the midpoint, and the height of the function at 1/3 and 2/3.

- Nondifferentiable function with an unlimited experimental budget, but where we want the fastest possible learning rate.

### 10.1.1 Differentiable functions

We begin by assuming that $f(x)$ is differentiable and unimodular (that is, it has a single local maximum). We assume the region we are searching is bounded, so we can scale the search region to be between 0 and 1. Initially, we assume the optimum $x^*$ can be anywhere in this interval with equal likelihood. If we measure the derivative $f'(x)$ at $x = 0.5$, we can observe whether $f'(.5) > 0$ or $f'(.5) < 0$. If the derivative is negative (as shown in Figure 10.1), then we know that $0 \leq x^* \leq .5$. We can eliminate the portion of the interval greater than .5. This means we can redefine the entire problem on the interval $(0, .5)$ and repeat the process. Let $\ell^n$ be the length of the interval in which the optimum may lie after $n$ iterations, where $\ell^0 = 1$. It is easy to see that $\ell^n = .5^n$. Under the assumptions of the problem, this is the fastest possible rate of reduction that we can achieve.

### 10.1.2 Nondifferentiable functions, finite budget

There are many problems where we cannot compute the derivative, but we can compute the function. For example, a transportation company may have a model that evaluates the on-time service when the fleet size is $x$. The company may vary $x$, re-running the model each time, but the model may not provide a derivative of the performance measure with respect to $x$. Instead, we have to search over a set of discrete values for $x$, just as we have done throughout this volume.

Normally we would create some sort of belief model, but now we are going to use the unimodular structure (along with our ability to perfectly observe the function). For example, we can try $x^1 = .2$ and $x^2 = .8$. If $f(.2) > f(.8)$, then we can eliminate the

region $(.8, 1)$ from further consideration. Of course, this only eliminates 20 percent of the interval. Our goal now is to eliminate the largest possible interval.

Imagine that we have an experimental budget of $N = 2$ experiments. In this case, the best strategy is to measure $x^1 = .5^-$ and $x^2 = .5^+$, by which we mean slightly less than .5 and slightly more than .5. Comparing $f(.5^-)$ and $f(.5^+)$ allows us to effectively eliminate half the interval, just as we did when we could compute a derivative (we are basically computing a numerical derivative).

Next consider the case where $N = 3$. If we first measure $x^1 = 1/3$ and $x^2 = 2/3$, we eliminate either $(0, 1/3)$ or $(2/3, 1)$. Assume that we eliminate the upper interval (as we did in Figure 10.1). Now we are left with the interval $(0, 2/3)$, but we have already run an experiment at the midpoint $x = 1/3$. We are going to use our final experiment at a point slightly above or below $1/3$ to determine the final interval of uncertainty, which will have width $1/3$. This last function evaluation will allow us to eliminate half of the remaining interval, giving us a region of width 1/6 where the optimum must lie.

We repeat this exercise one more time for $N = 4$, but now we are going to assume that the optimum is at $x = 0$ (but we have to construct our experiments without knowing this), so that we are always eliminating the upper part of the remaining interval. If we measure $x^1 = 2/5$ and $x^2 = 3/5$, we would eliminate $(3/5, 1)$, and we are left with the interval $(0, 3/5)$ with an experiment at $2/5$. Conveniently, this is at the two-thirds point of the interval $(0, 3/5)$, with two remaining experiments. If we measure $x = 1/5$, we are now in the same situation we were when we started with $N = 3$, but on an interval of width $3/5$. We eliminate the upper interval, leaving us with the interval $(0, 2/5)$ and an experiment at the midpoint $1/5$. We use our final experiment at a point slightly higher or lower than the experiment at $1/5$, giving us a final interval of $1/5$.

Now compare this result to what we would have obtained if we had used the bisection search with numerical derivatives. If $N = 2$, we would have measured just above and below .5, giving us a final interval of width .5. If $N = 4$, we would have done this twice, giving us an interval of width $.5^2 = .25$, which is greater than the interval we obtained of $1/5 = .2$. How did we accomplish this?

There is a pattern in this logic. When we eliminated the interval $(3/5, 1)$, we were left with the interval $(0, 3/5)$ and an experiment at $2/5$. Rescaling all the numbers so that the interval is of length 1, we get an interval of $(0, 1)$ with an experiment at $2/3$. If we eliminated the lower part of the interval $(0, 2/5)$, we would be left with the interval $(2/5, 1)$ (which still has length $3/5$) and an experiment at $3/5$. Rescaling gives us an interval of length 1, and an experiment at $1/3$. Either way, we end up with an experiment we would have made anyway if we only had 3 experiments.

We can formalize the algorithm as follows. Let $f^1 < f^2 < f^3 < \ldots$ be an increasing sequence of integers. If we are allowed three experiments, assume we first measure $f^1/f^3$ and $f^2/f^3$. If we eliminate the upper part of the interval, we are left with an interval $(0, f^2/f^3)$. If we eliminate the lower part of the interval, we are left with $(f^1/f^3, 1)$. We would like the width of these intervals to be the same, so we are

going to require that

$$\frac{f^2}{f^3} = 1 - \frac{f^1}{f^3},$$

or, rearranging slightly,

$$f^3 = f^1 + f^2.$$

Similarly if we have run four experiments, we first measure $f^2/f^4$ and $f^3/f^4$. Repeating the exercise above, rejecting the upper range leaves us with an interval of $f^3/f^4$, while rejecting the lower range leaves us with the interval $(1 - f^2/f^4)$. Again equating these gives us

$$\frac{f^3}{f^4} = 1 - \frac{f^2}{f^4},$$

or

$$f^4 = f^3 + f^2.$$

Using proof by extrapolation, we see that if we are making $N$ experiments, we want

$$f^N = f^{N-1} + f^{N-2}. \tag{10.1}$$

Furthermore, this has to be true for all $n < N$. Equation (10.2) defines what is known as the Fibonacci sequence, comprising the numbers $(1, 1, 2, 3, 5, 8, \ldots)$, where we initialize the sequence using $f^0 = 1$, and then let $f^1 = 1, f^2 = 2, f^3 = 3, \ldots$.

It is possible to show that the Fibonacci sequence produces an optimal search sequence for a finite experimental budget (under the assumptions of our problem). By this we mean that for a given value of $N$, no other method will produce a smaller final interval in which the optimal solution may lie. This means that our search is optimal, because we learn the most within our budget.

### 10.1.3 Nondifferentiable functions, infinite budget

So what if our budget is unlimited? We are going to start by hypothesizing that in the limit, we are going to measure the interval at two points, which we are going to denote by $1 - r$ and $r$, where $.5 < r < 1$. If we measure $1 - r$ and $r$, and then eliminate the upper interval, we are left with the interval $(0, r)$. Now assume we are going to measure the same points within this interval, which would occur at $r(1 - r)$ and $r^2$. We want the larger of these two points to coincide with the smaller of the two experiments in the original interval, so that at each iteration, we are measuring only one additional point (as we did with the Fibonacci search). This means that we require

$$r^2 = 1 - r,$$

or

$$r^2 + r - 1 = 0.$$

We solve this using the quadratic formula which gives the roots for the equation $ar^2 + br + c = 0$ as

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}.$$

Using only the positive part gives us the solution

$$
\begin{aligned}
r &= \frac{-1 + \sqrt{1 + 4}}{2} \\
&= \frac{-1 + \sqrt{5}}{2} \\
&= 0.618.
\end{aligned}
$$

The quantity $r = .618$ is known as the *golden section* or *golden ratio*. We note that we get the same result if we had eliminated the lower portion of the interval, since our experiments are naturally symmetric. We also note that if $f^n$ is the $n$th Fibonacci number, then

$$\lim_{n \to \infty} \frac{f^{n-1}}{f^n} \to r = 0.618.$$

Thus, the golden section search is the limit of the Fibonacci search.

The Fibonacci search and the golden section search are both examples of optimal learning in that they give the fastest improvement in our knowledge, measured by the length of the interval where we think the optimum might be. The Fibonacci search guarantees the smallest possible interval after a fixed (and known) number of experiments. The golden section search gives the fastest *rate* of convergence for an algorithm that will be run an infinite number of times. To put it another way, we account for the effect that our next experiment will have on the length of the uncertain interval.

## 10.2  NOISY EXPERIMENTS

We now return to our more familiar setting where experimental observations of the function are noisy. We consider the stochastic version of bisection search, where we assume we have access to a noisy indicator that tells us whether we think the optimum is to the left or the right of our experiment. In generalizing the bisection search we make a particular assumption about the experimental noise. We explain this assumption by noting that the bisection search operates by separating two regions of the search space: the region to the left of $x^*$, and the region to the right. Measuring the derivative of the function $f(x)$ at a point $x$ reveals in which part of the search space $x$ belongs.

Even without a function $f(x)$, we can use a bisection search to solve any problem where we want to find the boundary between two regions, and where measuring a point reveals in which region it resides. In the error-free case, this revelation is always correct. In the noisy version of the problem we instead assume that the revelation is incorrect with a probability that is *known* and *constant*. While these conditions are

typically not going to be satisfied in practice, they provide an elegant search model and they are certainly more realistic than assuming the observation of the function is perfect.

This assumption of constancy would tend to be met in applications where the transition between regions is abrupt. As an example, if a city's water supply were contaminated with a dangerous chemical we would want to localize the extent of contamination as quickly as possible, and if the chemical did not dissolve well in water but instead tended to stay concentrated, we would find a situation with this abrupt transition between contaminated and uncontaminated water. In contrast, when we measure a smooth function with additive noise, noise tends to cause incorrect region assignments more frequently near the function's maximum. With this in mind, we should be careful in applying the stochastic bisection algorithm presented here to situations not meeting its assumptions.

### 10.2.1   The model

We formulate our problem mathematically by again supposing that we have a point $x^*$ whose location is unknown beyond that it resides in the interval $[0, 1]$. The point $x^*$ corresponds to the boundary between the two regions in $[0, 1]$, so in the water contamination example, the water in region $(x^*, 1]$ would be contaminated and the water in region $[0, x^*]$ would not. We adopt a Bayesian prior density $p_0$ on the location of this point $x^*$, where $p_0(x)$ gives the likelihood (density) that $x^* = x$. It could be uniform on $[0, 1]$ if we had little real information about its location, or it could be some more complicated distribution expressing a stronger belief. We then suppose that we are offered the opportunity to take a sequence of experiments $x^0, x^1, \ldots, x^N$ in the interval, and with each experiment $x^n$ we get a noisy response $\hat{y}^{n+1}$ suggesting into which region $x^n$ resides. Given $x^n$ and $x^*$, this response will be independent of all other responses and will have the distribution

$$\hat{y}^{n+1} = \begin{cases} I_{\{x^* \leq x^n\}} & \text{with probability } q, \\ I_{\{x^* > x^n\}} & \text{with probability } 1 - q. \end{cases}$$

We may also express this as $\mathbb{P}\{\hat{y}^{n+1} = I_{x^* \leq x^n}\} = q$. Here $q$ is the probability our experiment is correct, and we assume this probability is known and constant. Equivalently, $1 - q$ is the error rate in our experiments.

### 10.2.2   Finding the posterior distribution

These experiments alter our prior belief about the location of $x^*$, giving us a posterior belief, all according to Bayes' rule. We use the notation $p^n$ to denote the posterior density at time $n$. To write the updating rule for $p^n$ explicitly, we introduce two pieces of notation. Let $F^n$ be the cumulative distribution function of the posterior at time $n$, by which we mean that, for any dummy variable $x$,

$$F^n(x) := \mathbb{P}\{x^* \leq x \mid p^n\} = \int_{[0,x]} p^n(x) \, dx.$$

We may also think of $F^n(x)$ as giving the probability that $x$ is in the region to the left of $x^*$, and in our water contamination example as giving the probability that the water at $x$ is not contaminated. As our second piece of notation, let $g$ be the function defined by $g(a, 1) = a$ and $g(a, 0) = 1 - a$, where $a$ will be a probability (such as the probability the true value is to the left or the right). Now we are ready to compute our updating rule.

Noting that nature's correct response to the experiment $x^n$ would be $I_{\{x^* \leq x^n\}}$, we write

$$\mathbb{P}\{\hat{y}^{n+1} = y \mid x^*, x^n\} = \begin{cases} q & \text{if } y = I_{\{x^* \leq x^n\}} \\ 1 - q & \text{if } y \neq I_{\{x^* \leq x^n\}} \end{cases}$$
$$= g(q, I_{\{y = I_{\{x^* \leq x^n\}}\}}).$$

So, $g(q, 1) = q$ corresponds to two cases: either $x^n < x*$, and we observe $\hat{y}^{n+1} = 1$ which correctly indicates that the optimum $x*$ is greater than our measured point; or $x^n > x*$, and $\hat{y}^{n+1} = 0$, which correctly indicates that the optimum $x^*$ is less than our measured point. The outcome $g(q, 0) = 1 - q$ corresponds to the opposite of both of these cases.

We may also write

$$\begin{aligned} \mathbb{P}\{\hat{y}^{n+1} = y \mid p^n\} &= \mathbb{P}\{x^* \leq x^n \mid p^n\} \mathbb{P}\{\hat{y}^{n+1} = y \mid x^* \leq x^n\} \\ &\quad + \mathbb{P}\{x^* > x^n \mid p^n\} \mathbb{P}\{\hat{y}^{n+1} = y \mid x^* > x^n\} \\ &= F^n(x^n) g(q, y) + (1 - F^n(x^n)) g(1 - q, y) \\ &= g(q F^n(x^n) + (1 - q)(1 - F^n(x^n)), y), \end{aligned}$$

where the last line may be seen by considering the cases $y = 0$ and $y = 1$ separately. Fixing some dummy variable $y$, we may then use Bayes' rule and these two relations to write,

$$\begin{aligned} p^{n+1}(x) \, dx &= \mathbb{P}\{x^* \in dx \mid p^n, \hat{y}^{n+1} = y\} \\ &= \frac{\mathbb{P}\{\hat{y}^{n+1} = y \mid p^n, x^* = x\} \mathbb{P}\{x^* \in dx \mid p^n\}}{\mathbb{P}\{\hat{y}^{n+1} = y \mid p^n\}} \\ &= \frac{g\left(q, I_{\{y = I_{\{x \leq x^n\}}\}}\right)}{g\left(q F^n(x^n) + (1 - q)(1 - F^n(x^n)), y\right)} p^n(x) \, dx. \end{aligned}$$

Substituting $\hat{y}^{n+1}$ for $y$ shows that our updating rule is

$$p^{n+1}(x) = \frac{g\left(q, I_{\{\hat{y}^{n+1} = I_{\{x \leq x^n\}}\}}\right)}{g\left(q F^n(x^n) + (1 - q)(1 - F^n(x^n)), \hat{y}^{n+1}\right)} p^n(x). \qquad (10.2)$$

The essential content of this updating rule is as follows. Our observation, if correct, would tell us into what region $x^n$ lies and would tell us whether $x^*$ is to the left or right of the experiment $x^n$. Let us give the name "suggested region" to the region in which the observation, if correct, would indicate $x^*$ resides. Since we know that the observation is only correct with probability $q$, we multiply the density in the
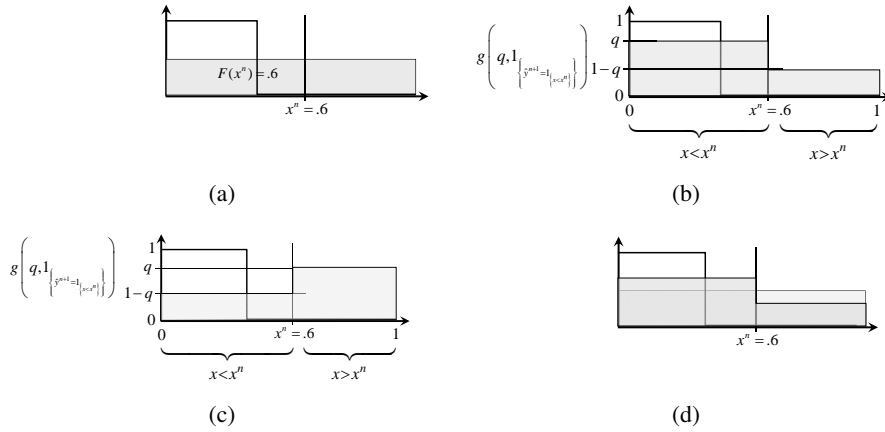
**Figure 10.2** Illustration of the effect of an observation of $\hat{y}$ for $x^n = .6$, starting with an initial uniform prior (a). (b) and (c) show the new beliefs for $\hat{y} = 1$ and $\hat{y} = 0$. The updated distribution $p^{n+1}(x)$ is given in (d).

suggested region by $q$ and the density in the other region by $1 - q$. This leaves us with a density which does not integrate to 1, so we then finish the update by normalizing.

As a numerical example, assume that $q = .7$, $x^n = .6$, and that we observe $\hat{y}^{n+1} = 1$. Assume that the distribution $p^n(x)$ is uniform, as depicted in Figure 10.2(a). Figures 10.2(b) and 10.2(c) show the conditional distribution given $\hat{y}^{n+1} = 1$ and $\hat{y}^{n+1} = 0$, respectively. Finally, Figure 10.2(d) shows the updated distribution for $p^{n+1}$ given $\hat{y}^{n+1} = 1$, overlaid on top of the original uniform distribution for $p^n(x)$. This distribution is computed using

$$
p^{n+1}(x) = \begin{cases} \frac{g(.7,1)}{g(.7(.6)+.3(.4),0)}p^n(x) = \frac{0.7}{.54}p^n(x) = 1.296 p^n(x) & x < x^n \\ \frac{g(.7,0)}{g(.7(.6)+.3(.4),0)}p^n(x) = \frac{0.3}{.54}p^n(x) = 0.556 p^n(x) & x > x^n \end{cases}.
$$

Let us briefly consider the case when our observations are always correct. Then $q = 1$ and we are back in the deterministic case. Suppose for the moment that $p^n$ is uniform on some interval $(a^n, b^n]$. Then our updating rule can be simplified to

$$
p^{n+1}(x) = \frac{g\left(1, I_{\left\{\hat{y}^{n+1}=I_{\{u \le x^n\}}\right\}}\right)}{g(F^n(x^n)), \hat{y}^{n+1})} \frac{I_{\{x \in [a^n, b^n]\}}}{b^n - a^n} = \begin{cases} \frac{I_{\{x \in (a^n, x^n]\}}}{x^n - a^n} & \text{if } \hat{y}^{n+1} = 1, \\ \frac{I_{\{x \in (x^n, b^n]\}}}{b^n - x^n} & \text{if } \hat{y}^{n+1} = 0.. \end{cases}
$$

Thus we see that the posterior is now still uniform but on some smaller interval, where either the points to the left or right of $x^n$ have been removed. Thus, if we begin with a uniform prior, i.e., with $p^0(x) = I_{x \in [0,1]}$, then our posterior at each time will again be uniform on a smaller interval. Comparing our knowledge in the deterministic case, where we had an interval in which we knew $x^*$ resided, to our knowledge here suggests that a natural way to express knowledge that $x^*$ lies in an interval is through a uniform probability distribution on that interval.

### 10.2.3  Choosing the experiment

Now let us return to the stochastic case where $q < 1$. With our updating rule in hand, we could take any given method for choosing our experiments, and compute the posterior density of the location of $x^*$ after our allotted number of experiments $N$, but before proceeding to say which experimental methods are best we need a way to evaluate the quality of the knowledge that we arrive to at the final time. In the deterministic case we knew that $x^*$ resided in an interval, and we evaluated how happy we were with our final knowledge according to the length of this interval. In the stochastic case, however, we no longer have an interval but instead a density expressing a continuum of belief about the location of $x^*$.

The objective function we use should correspond to our notion of length in the deterministic case, and it should punish greater uncertainty. There are many choices, but one possibility is the entropy. Denote the entropy of any particular density $p$ on the location of $x^*$ by $H(p)$,

$$H(p) := -\int_0^1 \log_2(p(x))p(x)dx$$

The entropy corresponds to uncertainty about $x^*$ in several senses. First, experiments always decrease entropy on average, no matter what is measured. Second, the entropy is largest for the uniform prior, which we may understand intuitively as the density we have when we are most uncertain about $x^*$. Third, the entropy approaches $-\infty$ as our posterior density sharpens to a point at the true location of $x^*$.

Additionally, the entropy corresponds in a very nice way to our use of interval length as the objective function in the deterministic case. In the deterministic case, we may know at a point in time that $x^*$ is in $[a, b]$, but we have no information about it beyond that. A natural belief to take on the location of $x^*$ in this situation is the uniform density on $[a, b]$, which is $p(x) = I_{x \in [a,b]}/(b - a)$. This density has entropy

$$H(p) = -\int_a^b \log_2(1/(b - a))/(b - a)dx = \log_2(b - a),$$

which is a strictly increasing function of the length $b - a$ of the interval. Thus minimizing the length of the interval $[b, a]$ is in some sense equivalent to minimizing the entropy. We characterize this equivalence more concretely later, when we show that the stochastic bisection algorithm is the same as the deterministic bisection algorithm when the probability $q$ of experimental error is $0$. Finally, an additional and very important reason for using entropy is that it provides an analytic and easy to use solution to the sequential problem.

Now with the transition function worked out and the entropy as our objective function, we have a well defined sequential information collection problem. This problem can be solved and the optimal solution computed using dynamic programming. Let us define our value function $V^n$ by taking $V^n(p^n)$ to be the smallest value of $\mathbb{E}\left[H(p^N) \mid p^n\right]$ than can be achieved starting from the density $p^n$ at time $N$. Bellman's principle then tells us that

$$V^n(p^n) = \min_x \mathbb{E}\left[V^{n+1}(p^{n+1}) \mid p^n, x^n = x\right].$$

In addition, we know that $V^N(p^N) = H(p^N)$ since there are no experiments left to make at time $N$.

We can use Bellman's recursion to compute $V^{N-1}$ from $V^N$. Since $V^N = H$, this recursion is $V^{N-1}(p^{N-1}) = \min_x \mathbb{E}\left[H(p^N) \mid x^{N-1} = x, p^{N-1}\right]$. Rather than computing this here, we simply state the following formula, which can be confirmed by direct computation.

$$\min_x \mathbb{E}\left[H(p^{n+1}) \mid x^n = x, p^n\right] = H(p^n) - q\log_2 q - (1-q)\log_2(1-q) - 1, \quad (10.3)$$

and the minimum is achieved by choosing $x$ to be a median of $p^n$. The median of $p^n$ is defined to be the point where $F^n(x) = 1/2$. If there is more than one median, any point satisfying $F^n(x) = 1/2$ achieves the minimum.

Now, using (10.3) and the Bellman recursion, we see that $V^{N-1}$ is given by

$$V^{N-1}(p^{N-1}) = H(p^{N-1}) - q\log_2 q - (1-q)\log_2(1-q) - 1,$$

and that the optimal decision $x^{N-1}$ is the median of $p^{N-1}$. Moreover, we see that the form of the value function at time $N-1$ is the same as it is at time $N$, but with a constant subtracted. This constant does not depend on $p^{N-1}$, nor does it depend on $n$. This tells us that if we repeat the computation of the Bellman recursion we will find that $V^{N-2}(p^{N-2}) = H(p^{N-2}) - 2\left(q\log_2 q + (1-q)\log_2(1-q) + 1\right)$, and that in general

$$V^n(p^n) = H(p^n) - (N-n)\left(q\log_2 q + (1-q)\log_2(1-q) + 1\right). \quad (10.4)$$

Furthermore, since the minimizer of the Bellman recursion at each time $n$ is the median of the density $p^n$ at that time, we have discovered that the optimal policy is to always measure at the median. Denoting the optimal experiment at time $n$ by $x^{*,n}$, we summarize this conclusion by saying that $x^{*,n}$ is such that

$$F^n(x^{*,n}) = 1/2.$$

Let us spend a few moments interpreting these results. First, the computation (10.3) tells us that if our goal is to minimize the expected entropy after one experiment, then the best $x^n$ has the expected posterior entropy equal to the original entropy at time $n$ minus a deterministic factor $q\log_2 q + (1-q)\log_2(1-q) + 1$. This factor is actually the mutual information between the experiment $\hat{y}^{n+1}$ and $x^*$, given that we measure at the median. This fact can be confirmed by computing the mutual information directly from its definition, although we do not perform this computation here.

We can view this reduction another way: the expected reduction in entropy about $x^*$ is equal to the information about $x^*$ contained in the outcome of the experiment, and this mutual information is maximized if we measure at the median of $x^*$. The mutual information is largest at the median because at this point we are "maximally uncertain" about to which region, $[0, x^*]$ or $(x^*, 1]$, it belongs since our belief assigns an equal probability of $1/2$ to each possibility. This conveys a general principle of information collection: often the experiment that is most valuable is the one whose result is least predictable. Put another way, if we already knew the result of an experiment before we took it, there would be no point in actually taking that experiment.

Then, the formula (10.4) shows that this general principle applies not just to single but multiple experiments. That is, the best we can do is to measure each time at the median of our belief, which is the experiment whose outcome is most uncertain, and the resulting decrease in expected entropy of $x^*$ is equal to the sum of the mutual information in all the experimental outcomes combined. From experiment $n$ onward, this decrease is $(N - n)\left(q \log_2 q + (1 - q) \log_2(1 - q) + 1\right)$ since there are $N - n$ experiments left to make and each contributes $q \log_2 q + (1 - q) \log_2(1 - q) + 1$.

We may also gain insight by taking the special case $q = 1$ and comparing to the deterministic case. As previously noted, when $q = 1$ and when we begin with a uniform prior on $[0, 1]$, the posterior remains uniform on a smaller interval. Then, the median of any such uniform posterior is simply the middle of the interval, and so we again always measure in the middle of the current interval, just as we did with deterministic bisection. Thus we see that the stochastic bisection algorithm reduces in the error-free case to exactly the classic bisection algorithm.

### 10.2.4 Discussion

We make one additional note about the usefulness of the stochastic bisection algorithm when our objective function is something different than the entropy objective we have assumed here. Although we have not shown it here, the decrease in entropy one obtains from measuring at the median is *deterministic*, even though the final density $p^N$ itself is certainly random. This is similar to the situation in the deterministic case, where the location of the final interval containing $x^*$ is unknown a priori, but the length of that interval is deterministic as long as we use the bisection rule.

This is a very nice property because it means that the optimal policy for the entropy objective function is also optimal for a broader class of objective functions. In particular, if our objective function is $\mathbb{E}[L(H(p^N))]$, where $L$ is some concave increasing function, then again one can show that the same stochastic bisection rule is optimal. We can think of $L$ as inducing some kind of risk aversion, in that using it would indicate we fear uncertainty about $x^*$ more than than we hope for certainty. For example, earlier we saw that the length of an interval is equal to the logarithm of the entropy of a uniform distribution on this interval, and so if we wanted to minimize the length of the final interval in the deterministic case, then perhaps we should minimize the logarithm of the entropy rather than just the entropy. But we have already minimized using this criterion, since the logarithm is concave and increasing and the stochastic bisection algorithm is optimal for this objective function as well.

## 10.3 BIBLIOGRAPHIC NOTES

Section 10.1 - For a proof of the optimality of the Fibonacci search sequence for unimodular functions, see Avriel & Wilde (1966).

Section 10.2 - The development in this section is due to Jedynak et al. (2011).

### PROBLEMS

**10.1**    How many iterations of a Fibonacci search sequence are needed to ensure that we can find the optimum of a unimodular function within 1 percent?

**10.2**    Consider the function $f(x|\alpha, \beta) = x^{\alpha-1}(1-x)^{\beta-1}$ for $0 \leq x \leq 1.0$. Perform the Fibonacci search to find the optimum of this function for the following values of $(\alpha, \beta)$.

   a) $\alpha = 8, \beta = 2$.

   b) $\alpha = 3, \beta = 12$.

   c) $\alpha = 6, \beta = 8$.

**10.3**    Repeat exercise 10.2 using the golden section search.

**10.4**    Assume that we are trying to find the maximum of the function $f(x) = x^6(1 - x)^2$, but now we are going to assume that we can compute the derivative.

   a) Using the derivative to indicate whether the optimum is to the left or the right, perform eight iterations of deterministic bisection search and report your best estimate of the optimal solution. Plot your distribution of belief describing where the optimum lies after eight observations.

   b) Now assume that we can estimate the sign of the derivative correction with probability $q = .70$ (even though we are computing it perfectly). Use the method described in Section 10.2 to find the optimum. Again plot your distribution of belief describing the location of the optimum.

   c) How would your answer to (b) change if $q = .5$? You should be able to answer this without repeating any of the calculations.