

# 1 Optimization

The field of *optimization* is large and vastly important, with a deep history in computer science (among other places). Generally, an optimization problem is defined by having a score function  $f : S \rightarrow T$  where  $S$  is some set of possible solutions to the problem and  $T$  is at least a partially ordered set of “scores.” Typically,  $T$  is a total ordering like the reals  $\mathbb{R}$  or the integers  $\mathbb{Z}$ . The goal of optimization is to find the globally best solution. Optimization algorithms are methods for sorting through  $S$  to find that best solution. There are several very general approaches to optimization, which differ mainly in (i) how they eliminate subsets of  $S$  from further exploration, (ii) how quickly they converge on a solution and (iii) how strong a guarantee they make on the quality of that solution.

The performance of an optimization technique depends, in large part, on the structure of  $f$ . When  $f$  is convex (Fig. 1a), the value of adjacent solutions is perfectly correlated, and greedy (“local”) approaches are always guaranteed to find the global optimum of  $f$ . (If a greedy algorithm is provably correct, then  $f$  must be convex.) However, most score functions are not convex, and the degree to which we can optimize our solution depends on how they are non-convex.

Consider the problem of finding the roots of some function  $f(x)$ . If  $f$  is continuously differentiable, then Newton’s method can be applied to construct a sequence of values  $x_0, x_1, \dots, x_n$  that converges on a root of  $f$ , i.e., converges on some  $x_*$  such that  $f(x_*) = 0$ .

To initiate this process, we must select some  $x_0$ . We then iterate the function

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)} , \tag{1}$$

to produce our sequence, which we halt when  $x_{n+1} - x_n = \Delta x < \epsilon$ , where  $\epsilon$  represents the desired accuracy in the root’s estimated value. But, Newton’s method is only guaranteed to find a root if certain conditions on  $f$  and  $x_0$  are fulfilled.<sup>1</sup> If we happen to choose  $x_0$  “sufficiently far” from a given root or if the derivative function is not “nice,” Newton’s method is not guaranteed to converge.

In general, a “nice” structure for  $f$  can be exploited to produce a fast algorithm for optimization, e.g., as in Newton’s method. In some cases, we may transform a non-convex score function  $g$  into a convex one  $f$ , e.g., in the max-flow/min-cut problem, where we transform the input graph  $G$  into the residual graph  $G_f$ , to which we can then apply a greedy optimizer. Max-flow/Min-cut is an

---

<sup>1</sup>The conditions are basically that the second derivative is defined everywhere, that  $x_0$  is “sufficiently close” (which has a technical definition) to the root  $r$ , and that the “basin of attraction” around  $r$  is well-defined (which also has a technical definition). To illustrate the idea of a basin of attraction, consider  $f(x) = x^3 - 2x^2 - 11x + 12$  and either  $x_0 = 2.352836$  or  $x_0 = 2.352837$ . These points lie in different basins, where a basin is defined as a set of initial choices that, when we apply Newton’s method, all converge on the same point.

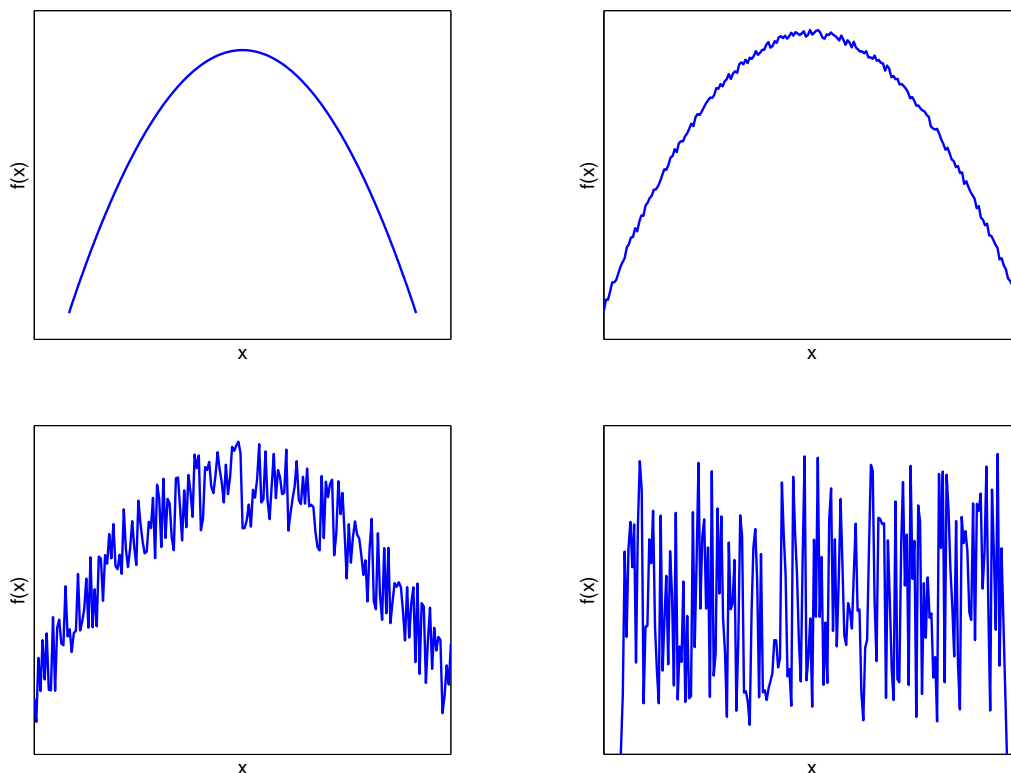


Figure 1: Examples of score functions with increasing “ruggedness.” The first shows a classic convex function, where the score of every point  $f(x)$  is tightly correlated with the score of an adjacent point, such that optimization algorithms based on local structure only (i.e., greedy optimizers) can always find the global optimum. Other figures show functions with progressively less correlation between the scores of adjacent points, ending in the last figure, in which adjacent points have uncorrelated scores.

important example in algorithms, as it demonstrates that thinking carefully about a problem can allow you to efficiently solve what would otherwise seem intractable. Another factor to consider in some cases is the computational time it takes to evaluate  $f(x)$ . If function evaluation is very expensive ( $\Omega(n^k)$  for  $k = \text{constant}$ , or, in practice, more than  $1/10000$  of the time you’re willing to wait for a full solution), then in addition to wanting a good solution, we want to find it while minimizing the number of function evaluations.

## 1.1 Exploration vs. Exploitation

Much of the qualitative behavior of optimization algorithms is well described by considering how they trade off “exploration” versus “exploitation.” Exhaustive-search algorithms, which evaluate  $f$  over the entire domain  $S$  and then pick the globally best solution, can be seen as maximally explorative, since they make no attempt to exploit any large-scale structure in  $f$  that would allow them to find the global optimum more quickly.

In contrast, “hill-climbing” or “gradient ascent” algorithms, which initialize at some point  $s \in S$  and then move greedily toward the nearest local optima, are maximally exploitative, since they make no attempt to explore alternative solutions away from its current trajectory.

Other techniques, including genetic algorithms, forward-backward algorithms, simulated annealing, the Nelder-Mead algorithm, etc. all strike some middle-ground between these two extremes. Generally speaking, algorithms that both explore and exploit, to some degree, can be applied to a wide variety of problems and expected to return reasonably good results.<sup>2</sup>

## 1.2 Simulated annealing

One fairly general optimization technique is *simulated annealing*.<sup>3</sup> The name comes from metallurgy: the “goal” in smithing is to produce a piece of metal (a sword, plow shears, a helmet, etc.) whose atomic configuration is very close to its minimum-energy structure; to get there, a smithy first heats the metal to a high temperature and then very slowly cools it, allowing the internal grain structures to align, thereby reducing the internal energy structure toward the global optimum.

The simulated version of this process works similarly, but with some slight differences because we generally are not using physics to define both the search algorithm and the score function, and instead must specify these for the problem at hand.

Specifically, we first require that we have some function  $f$  that can be called to return the score of an arbitrary input  $s$ . For simplicity, we assume calling  $f(s)$  is a constant time  $O(1)$  operation. We must also specify three additional components:

---

<sup>2</sup>Although I do not know of any completely general results here, it seems likely that a theorem should exist that relates the optimal tradeoff between exploitation and exploration and the large-scale smoothness of  $f$ . The boundary cases are clear: when  $f$  is convex, pure exploitation will always succeed, while when  $f$  is maximally rugged, with no correlation between the scores of adjacent points, pure exploration is the only way to succeed. The difficulty is rigorously handling the in-between structures on  $f$ .

<sup>3</sup>Simulated annealing can be derived directly from the Metropolis-Hastings algorithm used by the Manhattan project for simulating the behavior of nuclear piles, which was published by Metropolis in 1953.

1. a function `neighbor(s)` that takes as input some solution  $s$  and returns a “neighboring” solution<sup>4</sup>  $s'$ ,
2. a function `accept(f(s), f(s'), t)` that takes as input the scores of two states and a parameter  $t$  (the “temperature”) and makes a probabilistic decision about whether or not to accept the proposed move (i.e., to replace  $s$  with  $s'$ ), and
3. a function `decreaseTemperature()` that takes a real-valued parameter  $t$  and returns another real-valued  $t' < t$ . The monotonic structure of `decreaseTemperature()` is called the “cooling” or “annealing” schedule.

Crucially, for simulated annealing to work, we must require two behaviors, for all values of  $t$ , on the `accept` function:

$$\begin{aligned} \Pr(\text{accept} \mid f(s') > f(s)) &> \Pr(\text{accept} \mid f(s') < f(s)) \\ \Pr(\text{accept} \mid f(s') < f(s), t') &< \Pr(\text{accept} \mid f(s') < f(s), t) \text{ for } t' < t \end{aligned}$$

that is, `accept` tends to accept proposed states with better scores, and the probability of accepting a “bad” move decreases with time.

One of the most common choices for the `accept` function is the one used in the Metropolis algorithm for Markov chain Monte Carlo, which satisfies the *detailed balance* criterion for MCMC and maximizes acceptance probability, where

$$\begin{aligned} \Pr(\text{accept} \mid f(s') > f(s)) &= 1 \\ \Pr(\text{accept} \mid f(s') \leq f(s), t) &= e^{-|f(s')-f(s)|/t} \end{aligned}$$

that is, the algorithm always accepts good moves, and it accepts a bad move with probability that decreases exponentially with the difference in scores, normalized by the temperature  $t$ . Note that when  $t$  is large, the probability that we accept a move to a lower-scoring state will be close to 1 and thus the algorithm has no preference for better scoring states. In contrast, as  $t \rightarrow 0$ , the probability of accepting a move to a lower-scoring state approaches 0, and the algorithm reduces to a greedy hill-climbing strategy.

### 1.2.1 The algorithm

The idea is to use this search algorithm to construct a sequence of states  $s_0, s_1, \dots$  that converges on a final solution, much like Newton’s method except that here we construct this sequence probabilistically rather than deterministically. The choice of what state to add to our sequence is determined

---

<sup>4</sup>The choice of functions here has a large effect on both the time to convergence and on the probability that the algorithm converges on the global optimum.

by the `accept` function, which repeatedly considers replacing the current solution  $s$  with a “neighboring” solution  $s'$ . Finally, for the algorithm to exit, we must specify a convergence criterion, e.g.,  $f(s_{i+1}) - f(s_i) < \epsilon$ , where  $\epsilon$  is the desired numerical accuracy for the converged state. In general, a convergence criterion should be selected based on the particular application, and in some cases mathematically rigorous criteria exist. In most cases, a heuristic is chosen.

Here’s pseudocode for a generic simulating annealer:

```
Generic-SA(f,t0,s0) {
    (t,s) = (t0,s0)           // initial temperature, solution
    while not converged {
        s' = neighbor(s)      // propose new solution
        if accept(f(s),f(s'),t) { // accept?
            s = s'             // update solution
        }
        t = decreaseTemperature(t) // follow annealing schedule
    }
    return state
}
```

Initially, the probability of accepting a “bad” proposal is high and the algorithm’s behavior is close to an unbiased random walk on the solution space. Thus, the algorithm is highly explorative and the score of the current solution matters very little for determining the short-term trajectory of the optimization algorithm through the solution space. As the temperature decreases, so does the acceptance probability, and the algorithm becomes progressively more exploitative. In the limit of  $t \rightarrow 0$ , the probability of accepting a “bad” proposal approaches zero and the algorithm reduces to a hill-climbing approach that only changes states if the proposed state is strictly better. The derivative of the annealing schedule controls how quickly the algorithms transitions between exploration and exploitation.

The most common choice of cooling schedule is a geometric function, such that  $t_{i+1} = \alpha \cdot t_i$  where  $0 < \alpha < 1$ . When  $\alpha \approx 0$ , the temperature cools extremely quickly and the algorithm spends very little time exploring new solutions. When  $\alpha \approx 1$ , the schedule is more gentle and the algorithm explores more of  $S$ .

To many theorists, simulated annealing is horribly *ad hoc* because its behavior depends on uncontrolled choices for the annealing schedule, the acceptance function and the neighbor function. Poor choices of these can lead to dismal performance and bad solutions. For instance, if we consider the set of adjacencies on  $S$  induced by the neighbor function as a graph  $G$ , if  $G$  is not connected, i.e., if for at least one pair  $u, v \in S$  there is no  $u \rightsquigarrow v$ , then the choice of initial state  $s_0$  effectively determines the best solution possible, which is the best solution in the component containing  $s_0$ .

### 1.2.2 Hill climbing, aka gradient ascent

The classic optimization algorithm known as “hill climbing” or “gradient ascent” is simply a special case of simulated annealing with the Metropolis-Hastings acceptance function: let  $t = 0$ . Then, the probability of accepting a bad state is always zero and the algorithm only changes solutions if the proposed solution is better than the current solution.

It is not hard to prove that such a hill climbing algorithm is guaranteed to converge on a *local optimum* of  $f$  (do you see how?). More generally, this same convergence guarantee holds for simulated annealing under most cooling schedules. Some theoretical results exist showing that under an infinitely slow cooling schedule, where, if  $i$  indexes the number of steps of the algorithm,  $\partial t / \partial i \rightarrow 0^+$ , the algorithm converges on the global optimum. The convergence time, however, can be (infinitely) slow. For some classes of score functions  $f$ , these results can be improved; for instance, functions with fairly “smooth” surfaces tend to converge quickly and in some cases the convergence can be proved to be quite fast.

### 1.2.3 Variations

Variations on the basic simulated annealing algorithm are universally motivated by the desire to reduce the likelihood that the algorithm will get stuck in local optima. For instance, in *simulated tempering*, whose name is also an analogy with metallurgy, the requirement that the cooling schedule be monotonic is relaxed and instead we only require that as the number of algorithm steps  $i \rightarrow \infty$ , the temperature  $t \rightarrow 0$ . This allows us to choose a cooling schedule that periodically raises the temperature parameter  $t$ , which shifts the algorithm’s behavior back toward exploration, allowing it to escape local optima and their corresponding basins of attraction.

### 1.2.4 An example

As an example to demonstrate simulated annealing, suppose we want to find the global maximum of the function  $f(x) = \sin(x)/x$ . This is a function with a fairly rugged surface and a strong global maximum at  $x = 0$ .

Figure 1 shows the trajectory and results of a simple SA algorithm initialized at  $x_0 = 10$ . I used the Metropolis-Hastings acceptance function and the neighbor function  $x' = x + N(0, 1/25)$  where  $N(\mu, \sigma^2)$  is a Normal distribution with mean  $\mu$  and variance  $\sigma^2$ , and the cooling schedule  $t' = 0.99t$ . The red shows the surface of the score function, the blue circles show the trajectory of the algorithm and the single green circle shows the final state (my convergence criteria was  $t < 10^{-10}$ ). In this trial, the algorithm did a fairly good job of finding the global maximum, while on others, it would often get stuck in a local maximum. As it turns out, this is a particularly nasty choice of  $f$ . Do you see why?

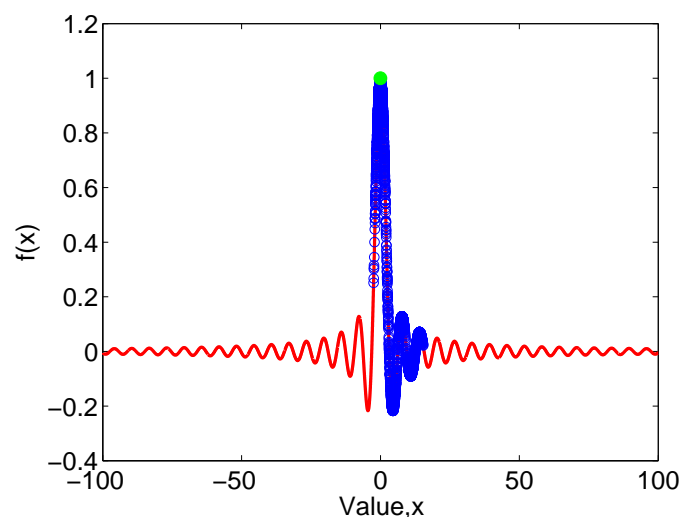


Figure 2: An example of running a simulated annealing algorithm to find the global maximum of the function  $f(x) = \sin(x)/x$  with  $x_0 = 10$ . The green dot shows the final state.

### 1.3 Genetic algorithms

*Genetic algorithms* are another popular optimization heuristic, inspired by the ability of natural selection to quickly find fairly good solutions to complex problems.<sup>5</sup>

These algorithms have very similar structure to the simulated annealing algorithm described above, but with a couple of crucial differences. First, instead of the algorithm tracking a single solution, a genetic algorithm tracks a population of solutions  $G_t$ . As with simulated annealing, a genetic algorithm proceeds in steps, indexed here by  $t$ ; at each step the score of each solution is calculated and then a new population of solutions is created. The creation of the new population is where the analogy with biology enters: in the classic version, the size of the next generation  $G_{t+1}$  is the same as  $G_t$ , and for each member of  $G_{t+1}$ , we select two “parent” solutions  $s, t$  from  $G_t$  each with probability proportional to its “relative fitness”  $f(s)/\sum_{s \in G_t} f(s)$ . (Many other variations exist.) The parents are first “crossed over”, in which some parts of  $s$  and the remaining parts of  $t$  are combined to produce a new “offspring”  $u$ . Before  $u$  is placed in  $G_{t+1}$ , each of its parts is “mutated” with some small probability.

The mutation operator thus plays the role of the neighbor function in simulated annealing, defining

---

<sup>5</sup>Genetic algorithms were originally developed by John Holland in 1975. For a good reference, see Melanie Mitchell’s *An Introduction to Genetic Algorithms*, MIT Press 1999.

the “local” neighborhood of solutions. In contrast, the crossover operator is very different: instead of taking a small step in  $S$ , crossover causes the new solution  $u$  to be some complicated and large jump across the space, at least with respect to the number of “mutation” operators that would be needed to yield an equivalent change in state. Thus, the mutation operator plays the role of encouraging exploitation, by generating minor variations on existing themes, while the crossover operator plays the role of encouraging exploration, by generating very large changes on the existing population.

The structure of the crossover operator has made the behavior of genetic algorithms difficult to analyze mathematically. Most analyses try to bound the overall behavior by letting the crossover rate go to zero. In this case, a genetic algorithm reduces to a kind of population-level simulated annealing at a fixed temperature. The behavior of this algorithm can be shown to be equivalent to those of the *replicator equation*, which have a close analogy with Bayesian learning, i.e., they converge exponentially on good solutions.

## 1.4 A panoply of methods

There are an enormous number of optimization algorithms available—too many to even attempt to list—and unfortunately this panoply is not very well organized from either a theoretical perspective or from a critical applied perspective. A good applied taxonomy might group methods by what kind of score functions  $f$  would cause them trouble, and a good theoretical taxonomy might group them by the kinds of guarantees they provide. Instead, techniques tend to be grouped by problem and field. In general, because optimization tasks are ones that algorithm designer often face, it is thus worthwhile to understand the particular techniques used on that particular class of problem, as these techniques often have nice properties for that context.<sup>6</sup>

For complicated score functions, simulated annealing and genetic algorithms, along with techniques we haven’t discussed like “ant colony optimization” and “particle swarm optimization” are often reasonable first choices, largely because they seem to work pretty well in practice.<sup>7</sup> As a word of advice, you should be suspicious of all such algorithms: without a theoretical basis for their success, it’s difficult to understand under what conditions we should expect them to behave badly, and, when they do behave badly, it’s not clear which of the many uncontrolled choices embedded in them are to blame.

---

<sup>6</sup>For instance, in machine learning, the Expectation-Maximization (EM) algorithm is a powerful and mathematically principled optimization technique, but it can only be applied to score functions representing likelihood functions with certain properties.

<sup>7</sup>Plus, they sound sophisticated.



## 2 For Next Time

1. Something fun