

Ant Colony Optimization and Stochastic Gradient Descent

Nicolas Meuleau and Marco Dorigo

IRIDIA, Université Libre de Bruxelles

Avenue Franklin Roosevelt 50, CP 194/6, B-1050 Brussels, Belgium

nmeuleau@iridia.ulb.ac.be, mdorigo@ulb.ac.be

Abstract

In this paper, we study the relationship between the two techniques known as ant colony optimization (ACO) and stochastic gradient descent. More precisely, we show that some empirical ACO algorithms approximate stochastic gradient descent in the space of pheromones, and we propose an implementation of stochastic gradient descent that belongs to the family of ACO algorithms. We then use this insight to explore the mutual contributions of the two techniques.

Keywords

Heuristic; Ant System; Ant colony optimization; Combinatorial optimization; Stochastic gradient descent; Reinforcement learning.

1 Introduction

The study of self-organization in social insects as a source of inspiration for novel distributed forms of computation is a promising area of AI known as *ant algorithms* (or sometimes as *swarm intelligence*) that is knowing growing popularity [4, 5, 9, 11]. A particularly successful form of *ant algorithms* are those inspired by the ant colonies foraging behavior. In these algorithms, applied to combinatorial optimization problems, a number of artificial ants are given a set of simple rules that take inspiration from the behavior of real ants. Artificial ants are then left free to move on an appropriate graph representation of the considered problem: they probabilistically build a solution to the problem and then deposit on the graph some artificial pheromones that will bias the probabilistic solution construction activity of the future ants. The amount of pheromone deposited and the way it is used to build solutions are such that the overall search process is biased towards the generation of approximate solutions of improving quality.

The historically first example of an algorithm inspired by the ants foraging behavior was ant system (AS) [8, 13] and its first application was to the Traveling Salesman Problem (TSP), a well known \mathcal{NP} -hard problem [19]. As a follow

up of AS, a number of similar algorithms, each one trying either to improve performance, or to make AS better fit a particular class of problems, were developed. Currently, many successful applications of such algorithms exist for \mathcal{NP} -hard academic combinatorial optimization problems such as quadratic assignment [16, 21], sequential ordering [14], resource-constrained project scheduling [22], vehicle routing with time windows [15], routing in packet-switched networks [7], shortest supersequence [23], and frequency assignment [20]. Applications to real-world combinatorial optimization problems are starting to appear: for example, a gasoline distribution company in Switzerland is using ACO algorithms to choose routes of its trucks [15], while Fujitsu-Siemens-Computers in Germany is testing ACO for an important logistic problem [26]. As a consequence, the ant colony optimization (ACO) metaheuristic was recently defined [10, 12] to put in a common framework all the algorithms that can be considered as offspring of AS.¹

Despite these successes, the basic mechanisms at work in ACO are still loosely understood, and there is usually no analytical tool to explain the observed effectiveness of an ACO algorithm. In this paper, we make a step in the direction of providing such tools by relating formally the ACO metaheuristic and the technique known as stochastic gradient descent (SGD), which has been extensively used in machine learning [24]. More precisely, we show that some empirical ACO algorithms approximate gradient descent of the expected value of the solution produced by an ant, in the space of pheromone trails. Moreover, we present an algorithm for combinatorial optimization that is, at the same time, a stochastic gradient descent algorithm working in the space of pheromones, and an instance of the ACO metaheuristic. This algorithm is an instance of the gradient-based reinforcement learning algorithm known as REINFORCE [32]. It can be seen both as a distributed, stigmergic² implementation of SGD, or as an ACO algorithm where the overall effect of ants activity is to descend the gradient of a given function in the space of pheromones.

The interest of establishing connections between ACO and SGD is that it offers many opportunities of cross-fertilization. On one side, many questions asked in the study of ACO algorithms receive a new light under the gradient-descent interpretation. For instance, a new way of understanding and proving convergence of ACO algorithms is proposed. Moreover, some classical acceleration techniques for gradient-based reinforcement learning can be easily transposed to ACO algorithms. On the other side, ACO algorithms show how to effectively implement the technique of SGD for solving large combinatorial optimization problems. Several improvement to the basic trial-and-error search of artificial ants developed in the best ACO algorithms suggest, in turn, new ways of using gradient descent in the framework of combinatorial optimization.

In this introductory paper, we present our main arguments using the example

¹A notable exception is the ABC algorithm for routing that, although belonging to ACO, was developed independently of AS [28].

²Stigmergy is a particular form of indirect communication used by social insects to coordinate their activities. Its role in the definition of ant algorithms is discussed at length in [4, 5, 9, 12].

of the asymmetric TSP. First (Section 2.1), we briefly review the ant system algorithm using the ATSP as example problem, and we show that AS is indeed closely related to the technique of SGD (Section 2.2 and 2.3). We then describe an implementation of SGD—or, more precisely, of William’s REINFORCE [32]—that belongs to the family of ACO algorithms (Section 2.4). In Section 2.5, we show how to generalize this reasoning to any combinatorial optimization problem that can be solved by an ACO algorithm. Finally, we comment these results and outline some future research directions in Section 3.

2 Ant System and Stochastic Gradient Descent

Ant system is a simple distributed algorithm that can be applied to any (constrained) minimum cost path problem on a graph. Throughout this paper, we use the application of AS to the asymmetric TSP (ATSP) as a basic example to present our arguments.

2.1 Ant System

The ATSP can be defined as follows. Let X be a set of cities, $|X| = n$, and $D = [d(x, y)]$ be a distance matrix, with $d(x, y) \in \mathbb{R}^+$ for all $(x, y) \in X^2$. We will denote by $\bar{X}^t \subseteq X^t$ the set of acyclic paths of length $t \in \{1, \dots, n\}$ in terms of the number of cities crossed ($\bar{X}^1 = X$ and $\bar{X}^2 = \{(x, y) \in X^2 : x \neq y\}$). ATSP can be defined as the problem of finding a path $h_n = (x_1, x_2, \dots, x_n) \in \bar{X}^n$ that minimizes the length of the corresponding tour, defined as

$$L(h_n) = \sum_{t=1}^{n-1} d(x_t, x_{t+1}) + d(x_n, x_1) .$$

The main variables of the AS algorithm are the pheromone trails $\tau(x, y)$ associated with each pair of cities $(x, y) \in \bar{X}^2$. Let \mathcal{T} be the bidimensional vector gathering all the $\tau(x, y)$ ’s. The basic principle of AS is to simulate artificial *ants* that use the pheromone trails to build a random tour. Once its tour completed, each ant makes a backward trip following the same path and updates the pheromones on its way back. Finally, the pheromone trails partially evaporate, that is, they decrease by a constant factor $\rho > 0$ called the *evaporation rate*. The behavior of each ant can be summarized as follows:

Forward:

- Draw the start city x_1 at random uniformly;
- At each step $t \in \{1, \dots, n-1\}$, after following the path $h_t = (x_1, x_2, \dots, x_t) \in \bar{X}^t$, draw the next city at random following

$$Pr(x_{t+1} = x \mid \mathcal{T}, h_t) = \begin{cases} 0 & \text{if } x \in h_t , \\ \tau(x_t, x) / \sum_{\substack{y \in X \\ y \notin h_t}} \tau(x_t, y) & \text{otherwise} , \end{cases} \quad (1)$$

where $x \in h_t$ means that the acyclic path h_t traverses x .

Backward:

- After generating the path $h_n = (x_1, x_2, \dots, x_n) \in \bar{X}^n$, reinforce the pheromone trails $\tau(x_t, x_{t+1})$ for each $t \in \{1, \dots, n-1\}$ and $\tau(x_n, x_1)$ by the amount $1/L(h_n)$.

There are many ways of implementing the algorithm. In the original implementation of AS, a set of m artificial ants synchronously built m solutions as follows: first, all the ants perform their forward trip without updating the pheromones, and then all of them execute their backward trip and update the pheromones for the next “generation” of m ants. A pheromone evaporation stage takes place at each generation, before sending the ants backward. The total update at each generation of each pheromone $\tau(x, y)$, $(x, y) \in \bar{X}^2$ is then

$$\Delta\tau(x, y) = \left(\sum_{i=1}^m \frac{\delta_{x,y}(h_n^i)}{L(h_n^i)} \right) - \rho\tau(x, y) ,$$

where $h_n^i \in \bar{X}^n$ is the path followed by the i^{th} ant during its forward trip ($i \in \{1, 2, \dots, m\}$), and $\delta_{x,y}(h_n) = 1$ if y is the immediate successor of x in the tour associated to $h_n \in \bar{X}^n$ and 0 otherwise. When $m = 1$, the different ants are sent one after the other in a fully sequential way, waiting for the previous ant to complete its backward trip before sending a new one. In this case, the pheromone update implemented by each ant is, for all $(x, y) \in \bar{X}^2$:

$$\Delta\tau(x, y) = \frac{\delta_{x,y}(h_n)}{L(h_n)} - \rho\tau(x, y) ,$$

where $h_n = \{x_1, x_2, \dots, x_n\} \in \bar{X}^n$ is the path followed by the ant during its forward trip. This pheromone update rule may also be used in a fully asynchronous and parallel implementation of ant system where ants act completely independently one on the other, and a pheromone evaporation stage is associated with each ant. It is often replaced by the following rule, that has proved more effective in practical applications [REF]:

$$\Delta\tau(x, y) = \rho \frac{\delta_{x,y}(h_n)}{L(h_n)} - \rho\tau(x, y) = \rho \left(\frac{\delta_{x,y}(h_n)}{L(h_n)} - \tau(x, y) \right) . \quad (2)$$

Note that the reinforcement of pheromone trails is now multiplied by the evaporation rate ρ .

Ant system was extensively tested together with a few other algorithms inspired by real ants behaviour on the TSP [13]. Although AS did not compete with the best known algorithms for TSP, its relative success inspired a great number of algorithms for different combinatorial optimization problems (cf. introduction). Often the AS-based algorithms provide state-of-the-art performance. These algorithms have recently been put in a unifying framework called ACO metaheuristic [10, 12]. ACO is composed of three main procedures.

In the first one, artificial ants probabilistically construct feasible solutions to the considered problem by moving on a proper graph representation. In this phase the construction process is biased by previous experience memorized in the form of pheromone trails, and, in some implementations, by heuristic information available about the considered problem (see discussion in section 3.1.1). The second phase, briefly discussed in section 3.1.2, is optional: here the solutions generated by the artificial ants can be taken to their local optima by a suitable local search routine. In the last phase, pheromone trails are updated by the ants, pheromone evaporation, and/or other suitable processes.

2.2 Stochastic Gradient Descent

ACO algorithms are usually regarded as optimization techniques working in the solution set of the combinatorial problem at hand. For instance, ant system for ATSP is usually seen as an algorithm that tries to find a tour of minimal length (that is, an optimal solution of the combinatorial problem). However, we adopt in this work a different approach and we look at ACO algorithms as working in the space of pheromone trails. In other words, we aim at finding *an optimal set of pheromones*, which can be defined in different ways. In this paper, we focus on a particular form of optimality for pheromone values. We will call an optimal set of pheromones a configuration that optimizes *the expected value of the solution produced by an ant during its forward trip*. We then study how this problem may be solved using gradient descent in the continuous space of pheromone trails.

In the case of ATSP, we aim at maximizing the expected value of the inverse length of an ant’s forward trip, given the current pheromone trails and the city-selection rule of equation (1). That is, we will *climb* stochastically the gradient of the “error” \mathcal{E} defined as:³

$$\mathcal{E} \stackrel{\text{def}}{=} \mathbb{E} \left[\frac{1}{L(h_n)} \mid \mathcal{T} \right] = \sum_{h_n \in \bar{X}^n} \Pr(h_n \mid \mathcal{T}) \frac{1}{L(h_n)} .$$

Note that the expectation is conditional on \mathcal{T} because the probability of a given tour happening depends on the current pheromone trail vector \mathcal{T} , while the “local error” $1/L(h_n)$ does not depend on the weights $\tau(x, y)$. Then we have

$$\frac{\partial \mathcal{E}}{\partial \tau(x, y)} = \sum_{h_n \in \bar{X}^n} \frac{\partial \Pr(h_n \mid \mathcal{T})}{\partial \tau(x, y)} \frac{1}{L(h_n)} ,$$

for each pair of cities $(x, y) \in \bar{X}^2$.

³Although we are dealing with maximization problems using stochastic gradient *ascent*, we use in this paper the vocabulary associated with gradient *descent* algorithms, which is much more common in the machine learning literature (see the example of artificial neural networks in section 3.2.1). The variable \mathcal{E} represents the objective function of the gradient ascent algorithm (and not of the original combinatorial optimization problem), and it must be maximized. However, it plays the same role as the error function used in gradient descent algorithms.

The probability of a given path is equal to the product of the probabilities of all the elementary events that compose it: if $h_n = (x_1, x_2, \dots, x_n) \in \bar{X}^n$, then

$$\Pr(h_n | \mathcal{T}) = \prod_{t=1}^n \Pr(x_t | \mathcal{T}, h_{t-1}) ,$$

where h_t is equal to h_n truncated after step t : $h_t = (x_1, x_2, \dots, x_t) \in \bar{X}^t$, and h_0 is the empty sequence. Therefore,

$$\frac{\partial \Pr(h_n | \mathcal{T})}{\partial \tau(x, y)} = \Pr(h_n | \mathcal{T}) \sum_{t=1}^n \frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} .$$

Here we have supposed that $\Pr(x_t | \mathcal{T}, h_{t-1}) > 0$, which is always true because $x_t \notin h_{t-1}$, as h_t is an acyclic path; and because the pheromone trails never fall to 0 in the original AS algorithm (however, we will see later that this is a problem for the new algorithm). Define the *eligibility trace*⁴ of (x, y) in path h_n as

$$T_{x,y}(h_n) \stackrel{\text{def}}{=} \frac{\partial \ln(\Pr(h_n | \mathcal{T}))}{\partial \tau(x, y)} = \sum_{t=1}^n \frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} , \quad (3)$$

then

$$\frac{\partial \mathcal{E}}{\partial \tau(x, y)} = \sum_{h_n \in \bar{X}^n} \Pr(h_n | \mathcal{T}) \frac{T_{x,y}(h_n)}{L(h_n)} = \mathbb{E} \left[\frac{T_{x,y}(h_n)}{L(h_n)} | \mathcal{T} \right] . \quad (4)$$

We will see later how to calculate the traces $T_{x,y}$. We can already outline the basis of the SGD algorithm. Climbing the gradient of \mathcal{E} corresponds to updating \mathcal{T} iteratively in the direction of the gradient of \mathcal{E} :

$$\Delta \mathcal{T} = \alpha \nabla_{\mathcal{T}} \mathcal{E} ,$$

that is,

$$\Delta \tau(x, y) = \alpha \frac{\partial \mathcal{E}}{\partial \tau(x, y)} ,$$

for each individual “weight” $\tau(x, y)$, where $\alpha > 0$ is the step-size parameter or *learning rate*. Following equation (4), we could do *exact* gradient ascent in the space of pheromone trails by enumerating all possible paths h_n and calculating, for each of them, the probability $\Pr(h_n | \mathcal{T})$ that an ant follows this path during its forward trip (given the current pheromone trails), the length of the corresponding tour $L(h_n)$, and the variable $T_{x,y}(h_n)$ for each $(x, y) \in \bar{X}^2$. Obviously, this approach would make no sense in practice: once we have enumerated all possible paths we can solve our original problem simply by taking the best.

⁴We borrow this vocabulary from reinforcement learning literature. It is used in similar gradient-based algorithm for optimal control Markov decision processes [1, 2, 32]. The *eligibility* of a weight is a measure of how much this weight will be involved in the next update. The $T_{x,y}$ variables are called *traces* because they keep track of the eligibility of the weights at each step t .

Moreover, the size of \bar{X}^n grows exponentially with the number of cities, so that this approach quickly becomes infeasible. Finally, exact gradient descent performs poorly in many complex domains because it gets trapped in the first local optimum on its way.

In *stochastic* gradient descent (SGD), an unbiased random estimate of the gradient is used instead of the true gradient. In the case of our application to ATSP, equation (4) shows that the gradient of \mathcal{E} is the expected value of random the random variable $T_{x,y}/L$ given the current pheromones (and the selection rule of equation (1)). Therefore, if we draw independently m paths $h_n^1, h_n^2, \dots, h_n^m$ in \bar{X}^n following the probability $\Pr(h_n \mid \mathcal{T})$, and average their *contributions* $R(h_n^i)T_{x,y}(h_n^i)$ to the gradient, then the result

$$\frac{1}{m} \sum_{i=1}^m \frac{T_{x,y}(h_n^i)}{L(h_n^i)}$$

is a random vector whose expected value is equal to the gradient. In other words, it is an unbiased estimate of the gradient. This is true regardless the number of paths sampled, even if only one sample is used to estimate the gradient (that is, $m = 1$). The resulting stochastic algorithm has a reasonable complexity.⁵ Moreover it may escape from some low-value local optima on its path. It sometimes makes bad moves because the gradient estimate is wrong, but these moves may allow jumping out of a bad local optimum. Therefore SGD performs usually better than exact gradient descent in large, highly multimodal, search spaces.

The basis of our comparison between ACO and SGD is the analogy between the actions of sending an ant forward and sampling a tour h_n from $\Pr(h_n \mid \mathcal{T})$. During its forward trip, the action of an ant is precisely to sample a solution following this probability distribution. Therefore, the forward component of AS can be used in a SGD algorithm as well, and we just have to change the weight update rules. We show below that the updates associated with a given sampled tour are very similar in the two algorithms.

2.3 A first ACO/SGD Algorithm

In our ACO/SGD algorithm, each artificial ant samples a tour $h_n \in \bar{X}^n$ using the current pheromones, and then updates every pheromone $\tau(x, y)$ following

$$\Delta\tau(x, y) = \alpha \frac{T_{x,y}(h_n)}{L(h_n)} \quad (5)$$

In a synchronous implementation, artificial ants are sent by groups of m ants that sample m tours without updating the pheromones (that is, following the same probability distribution $\Pr(h_n \mid \mathcal{T})$). Then, each of them update the

⁵We may actually influence the complexity by fixing the number m of samples drawn to calculate the estimate of the gradient. More samples allow a more accurate estimate.

pheromones as the function of the quality of the solution it generated during its forward trip. The total pheromone update rule is then

$$\Delta\tau(x, y) = \alpha \sum_{i=1}^m \frac{T_{x,y}(h_n^i)}{L(h_n^i)} ,$$

for all (x, y) , where $h_n^i \in \bar{X}^n$ is the path followed by the i^{th} ant. This corresponds to using m independent samples h_n to calculate the gradient estimate (note that the factor $1/m$ has been absorbed in the learning rate α). The particular case $m = 1$, that is based on equation (5), corresponds to a fully sequential implementation where only one sample is used to estimate the gradient and make a step in the space of pheromones. Equation (5) may also be used as the basis of a fully asynchronous and parallel implementation of the algorithm where the artificial ants act completely independently one on the other. In this case, the gradient estimates may be slightly biased by the simultaneous read and write activity of the different ants. However, this bias will probably be negligible in many applications, provided that the learning rate α stays in reasonable bounds.⁶

Given a path $h_n = (x_1, x_2, \dots, x_n) \in \bar{X}^n$ and a pair of cities $(x, y) \in \bar{X}^2$, the problem is now to calculate $T_{x,y}(h_n)$ as defined by equation (3). Using equation (1), we see that:

- if $x \neq x_{t-1}$ or $y \in h_{t-1}$, then $\Pr(x_t | \mathcal{T}, h_{t-1})$ is independent of $\tau(x, y)$ and

$$\frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} = 0 ;$$

- if $x = x_{t-1}$ and $y = x_t$ then

$$\begin{aligned} \frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} &= \frac{\partial \ln\left(\tau(x, y) / \sum_{y' \notin h_{t-1}} \tau(x, y')\right)}{\partial \tau(x, y)} , \\ &= \frac{\partial \ln(\tau(x, y))}{\partial \tau(x, y)} - \frac{\partial \ln\left(\sum_{y' \notin h_{t-1}} \tau(x, y')\right)}{\partial \tau(x, y)} , \\ &= \frac{1}{\tau(x, y)} - \frac{1}{\sum_{y' \notin h_{t-1}} \tau(x, y')} , \\ &= \frac{1 - \Pr(y | \mathcal{T}, h_{t-1})}{\tau(x, y)} ; \end{aligned}$$

- if $x = x_{t-1}$, $y \neq x_t$ and $y \notin h_{t-1}$ then, similarly,

$$\frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} = -\frac{\Pr(y | \mathcal{T}, h_{t-1})}{\tau(x, y)} ;$$

⁶Typically, the learning speed of a constant step-size SGD algorithm increases with the learning rate, up to a certain limit where the quality of the solution obtained is debased and the algorithm becomes unstable. By a “reasonable” learning rate value, we mean any value that is small enough to preserve the quality of the solution found and avoid instability.

Therefore, the SGD algorithm can be implemented using distributed, local information, as done in AS. The weight update corresponding to a sampled tour can be performed by the ant that sampled this tour during a backward trip. When returning to the t^{th} city of the tour⁷, the ant updates the pheromone trail $\tau(x_t, x)$ for all $x \in X$ following

$$\Delta\tau(x_t, x) = \alpha \frac{1}{L(h_n)} \frac{1}{\tau(x_t, x)} (\delta(x_{t+1} = x) - \Pr(x | \mathcal{T}, h_t)) \quad , \quad (6)$$

where $\delta(x_{t+1} = x) = 1$ if $x_{t+1} = x$ and 0 otherwise, and $\Pr(x | \mathcal{T}, h_t)$ represents the probability that the ant chooses $x \in X$ as the next city when it stands in x_t and it has already followed the path $h_t \in \bar{X}^t$. Note that the ants must be a little bit more “intelligent” than in the original AS. They need remember not only the tour they followed, but also the probability of choosing each city at each step of the forward trip. If they do not have such a memory, they can always recompute the probabilities using the pheromone trails, but the trails must not have changed in between due to another ant updating pheromones. Therefore, this solution works exactly only in a synchronous implementation of the algorithm.

The previous results show that SGD of the error $E[1/L(h_n) | \mathcal{T}]$ is close to the original AS algorithm. This appears clearly when we compare equation (6) to the analogous step-by-step update rule used in an asynchronous implementation of AS based on equation (2):

$$\Delta\tau(x_t, x) = \rho \left(\frac{\delta(x_{t+1} = x)}{L(h_n)} - \tau(x_t, x) \right) \quad . \quad (7)$$

The step-size α of SGD plays the role of the evaporation rate ρ in AS. The main differences are:⁸

- the pheromone value $\tau(x_t, x)$ in AS update rule is replaced in SGD by the probability $\Pr(x | \mathcal{T}, h_t)$ of moving from x_t to x ;
- the decrease of pheromones (through the term $-\tau(x_t, x)$ in (7), and through the term $-\Pr(x | \mathcal{T}, h_t)$ in (6)) is proportional to the *reward* $1/L(h_n)$ in our algorithm, while it is independent in the original AS;
- the presence of the factor $1/\tau(x_t, x)$ in the update rule (6) of the gradient-based algorithm.

It is important to note that, in our algorithm, the pheromones that are not used during an ant’s forward trip are not modified during the ant’s backward trip. If we had already visited y when we were in x , then $\tau(x, y)$ was not used to chose the next city after x , and hence, it was not used at all during the whole forward

⁷We assume here that the pheromone variables are stored as a set of tables $\mathcal{T}_x = [\tau(x, y)]_{y \in X}$, each \mathcal{T}_x being accessible to the artificial ants (or “physically located”) in city $x \in X$.

⁸Another difference between the two algorithms is that the pheromone $\tau(x_n, x_1)$ between the last and first city of the tour is reinforced by AS, while it is left untouched by SGD.

trip. As a consequence, $\tau(x, y)$ is left unchanged during the backward trip (that was not the case in the original AS where every pheromone trail evaporates). This makes sense, because if $\tau(x, y)$ is not used during the generation of a tour, then the value of this tour provides no information about the good way to move $\tau(x, y)$. This will be true in any application of SGD. It implies that the weight update associated with a forward trip (*i.e.*, a sampled solution), can always be performed in a backward trip following the same path. This is the basis of “ant” implementations of SGD presented in this paper.

There are a few problems with the algorithm we just defined. First, the evaporation rule may bring the weights $\tau(x, y)$ at or below 0. Negative pheromone trails do not really make sense. Moreover, we supposed the pheromones (strictly) positive when calculating the gradient. When some pheromones are 0, the analytical expression of the gradient is more complex. An empirical solution to this problem consists of preventing artificially the weights from falling below a given value $\epsilon > 0$. However, there is another problem with this algorithm: the contribution $T_{x,y}(h_n)/L(h_n)$ of a sequences $h_n \in \bar{X}^n$ may tend to infinity when some pheromone trails $\tau(x, y)$ tend to 0, which induces a very unstable behavior of the algorithm in some regions of the search space. For instance, if $\tau(x, y) \approx 0$ for some $(x, y) \in \bar{X}^2$ and an ant unluckily goes through this edge during its forward trip, then the subsequent weight update, that is proportional to $1/\tau(x, y)$, may be devastating and bring the algorithm very far from its original state. This is a case of unstable behavior due to unboundedness of the gradient estimate variance [3]: although the expected value of the gradient estimate (that is, the gradient itself) is always finite, the *variance* of the gradient estimate tends to infinity when some weights $\tau(x, y)$ tend to 0. In the next section, we present a new implementation of SGD that does not exhibit this instability, and that still belongs to the family of ACO algorithms.

2.4 A Stable ACO/SGD Algorithm

A classical solution to the problem of unbounded variance and unstable behavior is to use the Boltzmann law instead of the proportional selection rule of equation (1) (e.g., [1]). In our case, the city-selection rule takes the form:

$$\Pr(x_{t+1} = x \mid \mathcal{T}, h_t) = \begin{cases} 0 & \text{if } x \in h_t \text{ ,} \\ e^{\tau(x_t, x)} / \sum_{\substack{y \in X \\ y \notin h_t}} e^{\tau(x_t, y)} & \text{otherwise .} \end{cases} \quad (8)$$

The derivation presented in section 2.2 is still valid, the only changes are in the calculation of the traces (section 2.3):

- if $x \neq x_{t-1}$ or $y \in h_{t-1}$, then

$$\frac{\partial \ln(\Pr(x_t \mid \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} = 0 \text{ ;}$$

- if $x = x_{t-1}$ and $y = x_t$ then

$$\begin{aligned}
\frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} &= \frac{\partial \ln\left(e^{\tau(x, y)} / \sum_{y' \notin h_{t-1}} e^{\tau(x, y')}\right)}{\partial \tau(x, y)}, \\
&= \frac{\partial \ln(e^{\tau(x, y)})}{\partial \tau(x, y)} - \frac{\partial \ln\left(\sum_{y' \notin h_{t-1}} e^{\tau(x, y')}\right)}{\partial \tau(x, y)}, \\
&= 1 - \frac{e^{\tau(x, y)}}{\sum_{y' \notin h_{t-1}} e^{\tau(x, y')}} , \\
&= 1 - \Pr(y | \mathcal{T}, h_{t-1}) ;
\end{aligned}$$

- if $x = x_{t-1}$, $y \neq x_t$ and $y \notin h_{t-1}$ then, similarly,

$$\frac{\partial \ln(\Pr(x_t | \mathcal{T}, h_{t-1}))}{\partial \tau(x, y)} = -\Pr(y | \mathcal{T}, h_{t-1}) .$$

As in the previous case, the gradient-descent weight updates may be performed by the artificial ants during backward trips where they retrace their path backward. The new pheromone update rule is

$$\Delta\tau(x_t, x) = \alpha \frac{1}{L(h_n)} (\delta(x_{t+1} = x) - \Pr(x | \mathcal{T}, h_t)) . \quad (9)$$

It replaces equation (6). Note that the factor $1/\tau(x_t, x)$ has disappeared from the right hand side, making this rule very similar to the original AS update rule (equation (7)).

Ants need the same memory capacity as in the previous SGD algorithm. This new algorithm does not have the drawbacks of the previous. The weights $\tau(x, y)$ are unconstrained and can take any real value while keeping the probability of each path (strictly) positive. Moreover, the gradient estimate is uniformly bounded (that is, bounded by the same bound for all $h_n \in \bar{X}^n$) by the value $1/L^*$ where L^* is the length of the shortest path, and so, its variance is bounded by $1/L^{*2}$. Therefore, this algorithm is stable in any part of the search space.

Bounded variance is a necessary condition for convergence to a local optimum, but it is not sufficient [3]. As a matter of fact, it can be shown that stochastic gradient algorithms as ours may exhibit unbounded behavior, that is, some weights may tend to infinity [27]. A typical case of unbounded behavior is when the weights that generate an optimal solution tend to plus infinity, while all the other weights tend to minus infinity. The problem is that an optimal set of pheromone values (that is, a set that generates an optimal tour with probability 1) is obtained when the weights are positive or negative infinite. Therefore, the algorithm may continue to climb a pick infinitely, so that some weights diverge to infinity. However, this limitation has no consequence in our combinatorial optimization framework because we are not interested in having all ants follow an optimal path, but we just want that an optimal solution is generated (at least) once. Moreover, some classical tricks may be used to derive a formally convergent variant of our algorithm. Notably, we can:

- either artificially bound the weights away from infinity (that is, they are enforced to stay in a compact subset of \mathbb{R}^l , where l is the total number of weights), as in [27];
- or add to the objective function a penalty term in the form $-c \cdot \|\mathcal{T}\|$, where c is a constant and $\|\mathcal{T}\|$ is a norm of \mathcal{T} , so that the objective function tends to $-\infty$ when the weights tend to infinity [3].

We conjecture that any of these two solutions may be used to design a variant of our algorithm that converges with probability one to a local optimum of the error function.

2.5 Extensions

It is easy to modify the algorithm so that it optimizes other criteria than $E[1/L(h_n) \mid \mathcal{T}]$. For instance, if we want to minimize the expected tour length $E[L(h_n) \mid \mathcal{T}]$, then the update rule of the (stable) SGD algorithm becomes

$$\Delta\tau(x_t, x) = -\alpha L(h_n) (\delta(x_{t+1} = x) - \Pr(x \mid \mathcal{T}, h_t)) \ .$$

Here, the algorithm may be understood as maximizing the reward $-L(h_n)$, which is always negative.⁹ The shape of the objective function strongly determines the behavior of a gradient-following algorithm with constant step-size as ours. The performance of SGD may vary with different objective functions, even if these functions have the same local and global optima.

More generally, the same approach could be applied to every combinatorial optimization problem for which an ACO algorithm can be designed.¹⁰ The generic ACO/SGD approach to a given maximization problem with solution set \mathcal{S} and objective function $f : \mathcal{S} \rightarrow \mathbb{R}$ can be resumed as follows: First, design a stochastic controller that generates solutions in an iterative way using a set of weights \mathcal{T} . The controller is represented as a *construction graph* \mathcal{G} such that the generation of a solution corresponds to some sort of path in this graph, and the weights \mathcal{T} are attached to the arcs (or vertices) of \mathcal{G} . The weights, called pheromones, determine the transition probabilities in \mathcal{G} during the random generation of a solution. This first stage, which is called the choice of a *problem representation* in [12], is crucial. It transforms the static combinatorial problem $\max\{f(s) : s \in \mathcal{S}\}$ into the problem of optimizing a trajectory, that is, a dynamic problem.

The next step is to define the “error” function as the expected value of a solution produced by the controller, given the current weights:

$$\mathcal{E} = E[F(s) \mid \mathcal{T}] = \sum_{s \in \mathcal{S}} \Pr(s \mid \mathcal{T}) F(s)$$

⁹Conversely to the original AS, our gradient algorithm does not assume that the objective function is positive.

¹⁰That is, for which a constructive heuristic can be defined [10, 12]

where $F : \mathcal{S} \rightarrow \mathbb{R}$ strictly increases with f (that is, $f(s) > f(s') \implies F(s) > F(s')$).

Solutions s are generated by an iterative stochastic process that follows a finite number of steps. Let H be the set of trajectories in \mathcal{G} that the controller may follow when generating a solution, and $g : H \rightarrow \mathcal{S}$ be the function that assigns the solution produced to a given trajectory. The error may be rewritten as:

$$\mathcal{E} = \mathbb{E}[F(g(h)) \mid \mathcal{T}] = \sum_{h \in H} \Pr(h \mid \mathcal{T}) F(g(h)) .$$

The gradient of the error is then the expectation over all possible trajectories of the value of the solution produced multiplied by an eligibility trace T_τ :

$$\frac{\partial \mathcal{E}}{\partial \tau} = \sum_{h \in H} \Pr(h \mid \mathcal{T}) F(g(h)) T_\tau(h) = \mathbb{E}[F(g(h)) T_\tau(h) \mid \mathcal{T}] ,$$

for each individual weight τ . The trace $T_\tau(s)$ is the sum of the partial derivatives of the log of every elementary event that composes the trajectory s . Note that this decomposable structure of the gradient derives from the fact that solution generation is an iterative process, that is, from the very nature of the ACO approach.

Stochastic gradient descent can thus be implemented by sampling a few trajectories h —which can be seen as sending a few ants forward in the construction graph \mathcal{G} —and calculating their contribution to the gradient. In general, the weights that are not used when sampling a trajectory h have zero update according to the contribution of h . Therefore, the weight updates can be performed by artificial ants during backward trips in \mathcal{G} .¹¹ Finally, if the ants use the Boltzmann law to make random choices during forward trips, then the gradient estimate is uniformly bounded and the algorithm is stable in any part of the search space.

3 Discussion

Technically speaking, the ACO/SGD algorithm described above is not new. It is an instance of the generalized learning automaton [25] and, more precisely, of the gradient-based reinforcement learning algorithm known as REINFORCE [1, 32]. The originality of our work is to apply REINFORCE in the framework of ACO for combinatorial optimization, instead of the traditional Markov decision process (MDP) or partially-observable MDP (POMDP) used in the reinforcement learning literature [18, 30]. Therefore, this work establishes connections between

¹¹Note that the construction process does not necessarily have the same “taboo” aspect as in ATSP, where each weight is used at most once during an ant’s forward trip. The calculation of the SGD update is not more difficult in this case. The general rule is that if an ant traverses a vertex of \mathcal{G} (and uses the pheromone values stored here) several times during its forward trip, it must traverse this vertex (and update the pheromones stored here) the same number of times during its backward trip.

“classical” reinforcement learning and the less classical ACO learning. It suggests a general approach for applying reinforcement learning to combinatorial optimization which can be resumed as follows:

1. Design a parametrized stochastic controller that generates solutions *in an incremental way*, which turns the original static (optimization) problem into a dynamic (control) problem;
2. Use a reinforcement learning algorithm to (learn to) optimize the controller.

It is tempting to over-generalize the previous results and see in any ACO algorithm a more or less accurate approximation of the mechanism of gradient descent in the space of pheromones. In a sense, SGD is a very intuitive trial-and-error, Monte-Carlo technique that samples solutions, increases the probability of the best sampled solutions, and decreases the probability of the worst. Its particularity is to be grounded on a solid theoretical framework so that it is possible to give a sense to the update performed by the algorithm, but the basic intuition is the same as in ACO algorithms. Researchers in the field of ACO algorithms might find this position a little bit reductive. In fact, the best ACO algorithms are not limited to the simple trial-and-error ant search, but also use other optimization techniques such as constructive heuristics and local search routines [14, 29]. However, these features may also be grounded in the gradient-descent framework, and improve the algorithm performance (see Section 3.1). It is also interesting to note that the assimilation of ACO to approximate SGD allows us to draw a parallel with artificial neural networks (ANNs), because SGD is the basic principle behind the well-known backpropagation algorithm [6, 24]. Accordingly, we suggest in this work that the basic mechanisms at work in ACO and ANNs could be the same.

In the following, we discuss some opportunities of cross-fertilization between ACO and SGD, and then survey some important issues in the study of ACO algorithms.

3.1 Mutual Contributions

There are many opportunities of cross-fertilization between ACO and SGD (or, more precisely, gradient-based reinforcement learning). On one side, several efficient acceleration techniques for gradient-based reinforcement learning can easily be implemented in our ACO/SGD algorithm, and then generalized to other ACO algorithms (future work). On the other side, existing ACO algorithms suggest different ways of implementing the technique of SGD in the context of combinatorial optimization. The most successful applications of the metaheuristic are not limited to the simple trial-and-error ant search, but also use some “external” information in the form of constructive heuristics or (discrete) local search routines. They suggest different ways of merging SGD and these two combinatorial optimization techniques. In this section, we examine how our simple ACO/SGD algorithm can be used in combination with these techniques, as inspired by previous ACO algorithms.

3.1.1 Using Constructive Heuristics

Many successful implementations of the ACO metaheuristic combine information from the pheromone trails and heuristic information when generating solutions. In the case of AS, the city selection rule (1) is replaced by

$$\Pr(x_{t+1} = x \mid \mathcal{T}, \eta, h_t) = \begin{cases} 0 & \text{if } x \in h_t \text{ ,} \\ \tau(x_t, x)^a \eta(x_t, x)^b / \sum_{\substack{y \in X \\ y \notin h_t}} \tau(x_t, y)^a \eta(x_t, y)^b & \text{otherwise ,} \end{cases}$$

where $\eta \geq 0$ is a (constructive) heuristic function of $(x, y) \in \bar{X}^2$, and a and b are two (positive) parameters that determine the relative influence of pheromone trails and heuristic information. The function η reflects heuristic information about the good way to complete a partial solution. For instance, in the case of ATSP, a common choice is $\eta(x, y) = 1/d(x, y)$ for all $(x, y) \in \bar{X}^2$. In this case, the closest unvisited cities have larger probability to be chosen than without heuristic information. Moreover, in the successful applications of ACO to *non-stationary* (time-varying) problems, such as data packet routing in ANTNET [7], the function η is used to provide information to the algorithm about the current state of the problem.

There are several ways of integrating a similar mechanism in our algorithm. A particularly simple and elegant formulation is obtained when we replace the exponential selection rule (8) of our gradient algorithm by the following equation:

$$\Pr(x_{t+1} = x \mid \mathcal{T}, \eta, h_t) = \begin{cases} 0 & \text{if } x \in h_t \text{ ,} \\ e^{a\tau(x_t, x) + b\eta(x_t, x)} / \sum_{\substack{y \in X \\ y \notin h_t}} e^{a\tau(x_t, y) + b\eta(x_t, y)} & \text{otherwise ,} \end{cases} \quad (10)$$

where a and $b \geq 0$ are two external parameter that plays the same role as in the previous equation.¹² Equation (8) is obtained when $a = 1$ and $b = 0$. $b = 0$ and $a > 0$ corresponds to the Boltzmann law with a temperature parameter set to the value $1/a$ [1]. When $a = 0$ and $b > 0$, the algorithm does not use the pheromone trails at all. It is then a iterative stochastic heuristic search similar to the first stage of GRASP [REF]. Therefore, we dispose of a hole range of algorithms that extend from pure (heuristic free) gradient-based reinforcement learning, to simple (constant probability) stochastic heuristic search.

The next step is to recalculate the gradient to take into account the new selection rule. One more time, only the end of the calculation is modified. We have:

¹²Note that this formalism is compatible with negative heuristic functions. Therefore, we can chose $\eta(x, y) = -d(x, y)$ in the case of ATSP.

- if $x \neq x_{t-1}$ or $y \in h_{t-1}$, then

$$\frac{\partial \ln(\Pr(x_t | \mathcal{T}, \eta, h_{t-1}))}{\partial \tau(x, y)} = 0 ;$$

- if $x = x_{t-1}$ and $y = x_t$ then

$$\begin{aligned} & \frac{\partial \ln(\Pr(x_t | \mathcal{T}, \eta, h_{t-1}))}{\partial \tau(x, y)} \\ &= \frac{\partial \ln \left(e^{a\tau(x, y) + b\eta(x, y)} / \sum_{y' \notin h_{t-1}} e^{a\tau(x, y') + b\eta(x, y')} \right)}{\partial \tau(x, y)} , \\ &= \frac{\partial \ln(e^{a\tau(x, y) + b\eta(x, y)})}{\partial \tau(x, y)} - \frac{\partial \ln \left(\sum_{y' \notin h_{t-1}} e^{a\tau(x, y') + b\eta(x, y')} \right)}{\partial \tau(x, y)} , \\ &= a \left(1 - \frac{e^{a\tau(x, y) + b\eta(x, y)}}{\sum_{y' \notin h_{t-1}} e^{a\tau(x, y') + b\eta(x, y')}} \right) , \\ &= a(1 - \Pr(y | \mathcal{T}, \eta, h_{t-1})) ; \end{aligned}$$

- if $x = x_{t-1}$, $y \neq x_t$ and $y \notin h_{t-1}$ then, similarly,

$$\frac{\partial \ln(\Pr(x_t | \mathcal{T}, \eta, h_{t-1}))}{\partial \tau(x, y)} = -a \Pr(y | \mathcal{T}, \eta, h_{t-1}) .$$

This leads to the following update rule, after absorbing the constant factor a in the learning rate α :

$$\Delta \tau(x_t, x) = \alpha \frac{1}{L(h_n)} (\delta(x_{t+1} = x) - \Pr(x | \mathcal{T}, \eta, h_t))$$

The only difference with the update rule of the previous algorithm (equation (9)) is that the heuristic-independent probability $\Pr(x | \mathcal{T}, h_t)$ calculated following (8) is replaced by the heuristic-dependent probability $\Pr(x | \mathcal{T}, \eta, h_t)$ defined by (10). Therefore, the basic principle of the heuristic-free algorithm generalizes to the new selection rule: each ant has to memorize the probability distribution it uses at each step of its forward trip, and then decrease the pheromones on its way back proportionally to these distributions.

We have thus a simple framework that covers a great variety of algorithms. By varying the values of the external parameters a and b , we change two factors that strongly influence effectiveness of the algorithm. The first is the shape of the “error” function \mathcal{E} (defined as a function from pheromone vectors to real number) that the algorithm is approximately “descending”. It is not clear to us what is exactly the effect of the new update rule on the “landscape” we are exploring. A globally optimum is still obtained by putting infinite weights on the best paths (that do not change from one parameter setting to an other), but some important aspects such as the steepness of some peaks is modified.

The second and probably most important feature that is modified by varying parameters a and b is the sampling process used to estimate the gradient. The probability of sampling different paths and the update of the corresponding pheromones change from one setting to another. In the case of ATSP, the overall effect of using heuristic information (having $b > 0$) is that when an unvisited city y is close to the current city x , it has a bigger chance to be chosen than without heuristic information. In counterpart, the the pheromone trail $\tau(x, y)$ is decreased by a larger amount each time that y is considered as candidate successor to x . Further research is needed to understand how better performances may be obtained in this framework.

3.1.2 Using Discrete Local Search

The (ACO) metaheuristic is often used in conjunction with local search algorithms [14, 29]. In this approach, an ACO algorithm generates starting points for a discrete local search routine¹³. Each ant produces a solution, say s_1 , which is then transformed into another solution, say s_2 , by the local search. Then the pheromones are updated. As our goal is to maximize the quality of the final solution s_2 , pheromones updates must be proportional to the quality of s_2 , not s_1 . Given this, there are still two ways of updating the pheromones:

- either we reinforce the pheromones corresponding to the final solution s_2 . In other words, we do as if the solution s_2 was generated directly by the ant algorithm, without the help of the local search (in this approach, we suppose that there is a mapping between the solution set and the set of possible forward trajectories).
- or we reinforce the pheromones corresponding to the intermediate solution s_1 .

By analogy with similar procedures in the area of genetic algorithms [31], we call the first alternative the Lamarckian approach, and the second the Darwinian approach.

The main argument supporting the Lamarckian approach is that it is reasonable to think that, if the ant algorithm can be trained directly using the better solution s_2 , then it would be stupid to train it using the worse solution s_1 . In fact, in published ACO implementations, only the Lamarckian alternative has been used. In the case of SGD, however, the Darwinian approach may make more sense. It is easy to show that, if we try to maximize the expected value of the solution s_2 produced by the local search algorithm, then the update rule of an SGD algorithm is to reinforce the pheromones corresponding to the intermediate solution s_1 proportionally to the value of the final solution s_2 . The formal framework developed in section 2.5 can be used for this calculation, the effect of the local search being modeled in the function F . Having understood this,

¹³Gradient descent is itself a local search procedure, but it operates in the continuous space of controller weights, while the discrete local search used here operates in the discrete set of solutions.

we can derive qualitative arguments in favor of the Darwinian approach. For instance, if the good starting points of the local search are very far from the corresponding local optima in the topology of the gradient algorithm, then the Darwinian approach could outperform the Lamarckian.

3.2 Important Issues

By establishing connections with ANNs on a side, and reinforcement learning on another side, we also show that ACO algorithms are concerned with two important issues paradigmatically illustrated in these techniques. They are, respectively, the issue of *generalization* and *the exploration vs. exploitation dilemma*. In this section, we examine how these problems arise in ACO algorithms. It is clear that *any* reinforcement learning algorithm for combinatorial optimization has to deal with these two issues simultaneously.

3.2.1 Generalization

The most famous application of SGD is surely the algorithm known as backpropagation in artificial neural networks (ANNs), and the issue the most studied in backpropagation is probably the ability of the algorithm to generalize over inputs [6]. Stated simply, the problem is to find a set of weights $W = [w]$ for a network encoding a function F from an input set I to an output set O (say, $O = \mathbb{R}$), so that F approximates as much as possible a target function $F^* : I \rightarrow O$. Backpropagation learns an optimal configuration of weights by observing a set of training examples, that is, pairs $(i, F^*(i))$ with $i \in I$, and memorizing and generalizing these observations. In general, the input set I is a huge combinatorial set, if not an infinite set. Therefore, it is not possible to present every instance $i \in I$ in the training set. However, backpropagation is able to generalize the observed data to unseen instances. That is, it assumes that any unseen input i has a value $F^*(i)$ that is close to the value of the observed examples which are similar to i in some sense. It is well known that the ability of an ANN to generalize and its efficiency in generalization strongly depend on the network structure [6].

As we said, the backpropagation algorithm is an instance of SGD. More precisely, its overall effect is to descend the gradient of the mean square error

$$\mathcal{E}_{MS} = \sum_{i \in I} p_i (F(i) - F^*(i))^2 = \mathbb{E} [(F - F^*)^2 | W] \quad ,$$

where $(p_i)_{i \in I}$ is a given probability distribution on instances ($\sum_{i \in I} p_i = 1$). Note that, conversely to the case of our ACO/SGD algorithm, the expectation in this equation is conditional on the weights W because the local error

$$e_{MS} \stackrel{\text{def}}{=} (F - F^*)^2$$

depends on the weights, while the probabilities p_i do not. Therefore,

$$\frac{\partial \mathcal{E}_{MS}}{\partial w} = \sum_{i \in I} p_i \frac{\partial e_{MS}(i)}{\partial w} = \mathbb{E} \left[\frac{\partial e_{MS}}{\partial w} | W \right] \quad .$$

This result suggests an exact gradient algorithm that enumerates all possible inputs $i \in I$ for each step of gradient descent. Conversely, backpropagation uses an unbiased estimate of the gradient obtained by sampling a unique input $i \in I$. After sampling input i and comparing the actual output $F(i)$ and the desired output $F^*(i)$, backpropagation updates the weights of the network following

$$\Delta w = -\alpha \frac{\partial e_{MS}(i)}{\partial w} .$$

for each weight w . The values of $F(i)$ and $F^*(i)$ are used here to calculate the partial derivative $\partial e_{MS}(i)/\partial w$.

We see that the basic principles of backpropagation and our algorithm are the same. It is, in both case, a Monte-Carlo estimation of the gradient of a given error function, with respect to a set of weights attached to the components of a graph (the construction graph \mathcal{G} in one case, and the ANN itself in the other case). Also, both are distributed and parallel implementation of this principle. Therefore, our algorithm should have, at least in part, the same ability to generalize observed data over unseen instances as backpropagation.

Actually, it is not hard to convince oneself that generalization is as big an issue in ACO in general as in ANNs. For instance, it is clear that the application of AS to the ATSP works by generalizing the observed solutions: if a majority of the sampled tours that traverse a given arc $(x, y) \in \bar{X}^2$ are of good quality, then the algorithm increases the probability of traversing this arc. In a sense, it assumes that, *in general*, the tours that traverse (x, y) are of good quality. That is, it generalizes the observations.

As for backpropagation, the ability and efficiency of an ACO algorithm to generalize is mostly determined by the structure of the graph, that is, the problem representation used by the artificial ants. Moreover, we believe that it is usually easier to predict the kind of generalization that an ACO algorithm can do given its problem representation (as we just did for the ATSP), than it is to predict the effect of a given neural network architecture for a given problem.

3.2.2 Exploration vs. Exploitation

As we stressed above, the ACO metaheuristic can be seen as a way of applying reinforcement learning to combinatorial optimization problems. Thus, every ACO algorithm has to deal with one of the main issues in reinforcement learning: the exploration vs. exploitation dilemma [17]. This is the problem of finding an optimal compromise between the necessity to experiment the least known solutions (exploration), and the necessity to maximize rewards by taking the estimated best actions (exploitation). This problem is characteristic of real-time on-line learning, where one motivation is to learn as fast as possible, and another is to maximize the reward received during each experience.

In the case of ACO algorithms, the mechanisms of the exploration vs. exploitation dilemma are intimately linked with those of generalization. For instance, in the application to the ATSP, since the pheromones are attached to pairs of cities $(x, y) \in \bar{X}^2$, the algorithm is confronted to questions like: “Is it necessary to

try again a given pair of cities that seems to be non-optimal?”. Clearly, other questions will arise with another problem representation. It appears that an optimal solution to the exploration vs. exploitation dilemma in the framework of ACO depends intimately on the problem representation, that is, the structure of the construction graph.

In practice, our ACO/SGD algorithm does not really address the issue of exploration, although it does face the dilemma. Its behavior is dictated by the trial-and-error search of SGD, independently of any consideration about exploration.

4 Conclusions

In this paper, we explored the connections between the two techniques of ACO and SGD. First, we showed that the empirical designed AS algorithm is very similar to SGD in the space of pheromones, and we proposed a stable implementation of gradient-based reinforcement learning that belongs to the framework of ACO algorithms. Then we outlined a general ACO/SGD algorithm for combinatorial optimization. The performance of this algorithm depends crucially on some basic choices such as the problem representation and the objective function. This insight may be used to develop simple acceleration techniques for ACO algorithms, by transposing previous work on gradient-based reinforcement learning. Moreover, the most successful applications of the ACO metaheuristic suggest new ways of merging gradient descent with other optimization techniques for combinatorial optimization.

In conclusion, we believe that this work constitutes a significant step toward understanding the mechanisms at work in ACO algorithms, bringing a new light to some important issues in the theory of these algorithms.

Acknowledgments

This work was supported by a Marie Curie Fellowship awarded to Nicolas Meuleau (Contract No. HPMFCT-2000-00230). Marco Dorigo acknowledges support from the Belgian FNRS, of which he is a Senior Research Associate. This work was also partially supported by the “Metaheuristics Network”, a Research Training Network funded by the Improving Human Potential programme of the CEC, grant HPRN-CT-1999-00106. The information provided is the sole responsibility of the authors and does not reflect the Community’s opinion. The Community is not responsible for any use that might be made of data appearing in this publication.

References

- [1] L. C. Baird and A. W. Moore. Gradient descent for general reinforcement learning. In *Advances in Neural Information Processing Systems, 11*, pages

968–974. MIT Press, Cambridge, MA, 1999.

- [2] J. Baxter and P. L. Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15:319–350, 2001.
- [3] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, MA, 1995.
- [4] E. Bonabeau, M. Dorigo, and G. Theraulaz. *Swarm Intelligence: From Natural to Artificial Systems*. Oxford University Press, New York, NJ, 1999.
- [5] E. Bonabeau, M. Dorigo, and G. Theraulaz. Inspiration for optimization from social insect behavior. *Nature*, 406:39–42, 2000.
- [6] K. Boyanov, editor. *Parallel and distributed processing*. North-Holland, Amsterdam; New York, 1990.
- [7] G. Di Caro and M. Dorigo. AntNet: Distributed stigmergetic control for communications networks. *Journal of Artificial Intelligence Research*, 9:317–365, 1998.
- [8] M. Dorigo. *Optimization, Learning and Natural Algorithms* (in Italian). PhD thesis, Dipartimento di Elettronica, Politecnico di Milano, Italy, 1992. pp. 140.
- [9] M. Dorigo, E. Bonabeau, and G. Theraulaz. Ant algorithms and stigmergy. *Future Generation Computer Systems*, 16(8):851–871, 2000.
- [10] M. Dorigo and G. Di Caro. The Ant Colony Optimization meta-heuristic. In D. Corne, M. Dorigo, and F. Glover, editors, *New Ideas in Optimization*, pages 11–32. McGraw Hill, London, UK, 1999.
- [11] M. Dorigo, G. Di Caro, and T. Stützle (Editors). Special issue on “Ant Algorithms”. *Future Generation Computer Systems*, 16(8), 2000.
- [12] M. Dorigo, G. Di Caro, and L. M. Gambardella. Ant algorithms for discrete optimization. *Artificial Life*, 5(2):137–172, 1999.
- [13] M. Dorigo, V. Maniezzo, and A. Colorni. The Ant System: Optimization by a colony of cooperating agents. *IEEE Transactions on Systems, Man, and Cybernetics – Part B*, 26(1):29–41, 1996.
- [14] L. M. Gambardella and M. Dorigo. Ant Colony System hybridized with a new local search for the sequential ordering problem. *INFORMS Journal on Computing*, 12(3):237–255, 2000.
- [15] L. M. Gambardella, È. D. Taillard, and G. Agazzi. MACS-VRPTW: A multiple ant colony system for vehicle routing problems with time windows. In D. Corne, M. Dorigo, and F. Glover, editors, *New Ideas in Optimization*, pages 63–76. McGraw Hill, London, UK, 1999.

- [16] L. M. Gambardella, È. D. Taillard, and M. Dorigo. Ant colonies for the quadratic assignment problem. *Journal of the Operational Research Society*, 50(2):167–176, 1999.
- [17] L. P. Kaelbling. *Learning in Embedded Systems*. MIT Press, Cambridge, MA, 1993.
- [18] L. P. Kaelbling, M. L. Littman, and A. W. Moore. Reinforcement learning: a survey. *Journal of Artificial Intelligence Research*, 4:237–285, 1996.
- [19] E. L. Lawler, J. K. Lenstra, A. H. G. Rinnooy Kan, and D. B. Shmoys. *The Travelling Salesman Problem*. John Wiley & Sons, Chichester, UK, 1985.
- [20] V. Maniezzo and A. Carbonaro. An ANTS heuristic for the frequency assignment problem. *Future Generation Computer Systems*, 16(8):927–935, 2000.
- [21] V. Maniezzo and A. Colorni. The Ant System applied to the quadratic assignment problem. *IEEE Transactions on Data and Knowledge Engineering*, 11(5):769–778, 1999.
- [22] D. Merkle, M. Middendorf, and H. Schmeck. Ant colony optimization for resource-constrained project scheduling. In *Proceedings of the Genetic and Evolutionary Computation Conference (GECCO-2000)*, pages 893–900. Morgan Kaufmann Publishers, San Francisco, CA, 2000.
- [23] R. Michel and M. Middendorf. An ACO algorithm for the shortest super-sequence problem. In D. Corne, M. Dorigo, and F. Glover, editors, *New Ideas in Optimization*, pages 51–61. McGraw Hill, London, UK, 1999.
- [24] T. Mitchell. *Machine Learning*. McGraw-Hill, Boston, MA, 1997.
- [25] K. S. Narendra and M. A. L. Thathachar. *Learning Automata: An Introduction*. Practice Hall, Englewood Cliffs, NJ, 1989.
- [26] R. Palm. Personal communication. 2002.
- [27] V. V. Phansalkar and M. A. L. Thathachar. Local and global optimization algorithms for generalized learning automata. *Neural Computation*, 7:950–973, 1995.
- [28] R. Schoonderwoerd, O. Holland, J. Bruten, and L. Rothkrantz. Ant-based load balancing in telecommunications networks. *Adaptive Behavior*, 5(2):169–207, 1996.
- [29] T. Stützle and H. H. Hoos. The *MAX-MIN* Ant System and local search for the traveling salesman problem. In *Proceedings of the 1997 IEEE International Conference on Evolutionary Computation (ICEC'97)*, pages 309–314. IEEE Press, Piscataway, NJ, 1997.

- [30] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, 1998.
- [31] D. Whitley, S. Gordon, and K. Mathias. Lamarckian evolution, the Baldwin effect and function optimization. In *Proceedings of PPSN-III, Third International Conference on Parallel Problem Solving from Nature*, pages 6–15. Springer Verlag, Berlin, Germany, 1994.
- [32] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8(3):229–256, 1992.