

Getting Reinforcement Learning to Work on Real Robots

William D. Smart Leslie Pack Kaelbling
Computer Science Department
Brown University
Providence, RI 02912
{wds, lpk}@cs.brown.edu

Submission to Workshop W7: Robot Exploration and Learning

1 Introduction

Robots are now begin used in complex, unstructured environments, performing ever more sophisticated tasks. As the task and environmental complexity increases, the need for learning on such systems is becoming more and more apparent. Robot programmers often find it difficult to articulate their knowledge of how to act in a form suitable for robots to use. Even when they can, the limitations on robot sensors and actuators might render their intuitions useless. Thus, having the robot be able to learn to act seems like a good idea.

One promising approach to learning on robots is *reinforcement learning* (RL). All that the programmer has to supply is a reward function, which maps states of the world onto a scalar reward value. Designing such a function seems more intuitive than writing code to perform a task, since we have a reasonably good idea of what is “good” and “bad” for the robot even if it is hard to explicitly say what it should be doing in detail. For example, in a particular task, it might be good to press a certain button (meaning a positive reward), quite bad to sit still doing nothing (a small negative reward) and *very* bad to fall down the stairs (a large negative reward). Armed with this reward function, we can use trial-and-error methods to learn a mapping from states of the world to robot actions with the intention of maximizing the total reward that we earn.

Although learning is useful for robots, it is not appropriate in all situations. For the work described in this paper, we are interested in improving the performance of the robot on some specific task. With this in mind, the use of learning is only appropriate in the following situations.

- When the task is originally under-specified or is difficult to code exactly by hand.
- When the task has parameters that are likely to change over time in unpredictable ways.
- When the time taken to learn the control policy is less than that taken to hand-code a comparable policy.
- When the learned policy can be executed more efficiently than a hand-coded one.

The last point requires some clarification. We may be able to conveniently write a control policy using, for example, a large rule-based system to generate actions. If we could learn the policy that these rules embody in a simple neural network, then we might prefer the learned policy because, once learned, it could generate actions faster than the rule-based system could.

There are several problems in using RL on robots, however. In the next section, we outline some of these difficulties. We then go on to describe a framework addressing these problems and conclude by discussing some of the issues that must be resolved before we can effectively use reinforcement learning in this framework.

2 Reinforcement Learning on Robots

As stated above, reinforcement learning uses a reward function to guide a trial-and-error search for the best action to take from each possible state of the world. Although systems utilizing RL have met with considerable success in many problem domains, there are some serious difficulties to be overcome when dealing with RL on robots.

In order to determine what the best action in a given state is, the learning system must have tried each possible action from that state at least once. On a robot, the number of states is typically the cross product of all individual sensor states. In general, there are many sensors, each with many states, leading to a huge number of possible states for the robot to be in. In addition, there are often many actions that the robot can take from each of these states. For example, the robot which we use has 24 SONAR sensors (in addition to several others). If we quantize the readings from these into 3 range bands (close, medium and far), then this gives a total of over 282 billion possible states that the robot can be in. From each of these states, if we can go either forwards or backwards, this leads to over 560 billion state-action combinations to try.

Since the robot is a physical system, actions take a certain time to perform. Assuming that we can perform one action each second, to try each of these possibilities would take somewhere in excess of 20,000 years. Obviously we cannot hope to visit each of these possibilities. In reality, we do not have to visit all of the states, as long as we assume certain continuities in the mappings we are trying to learn. However, the number of state-action pairs to be tried even under these assumptions often proves to be too many to accomplish in practice.

A related problem is seen in early stages of learning. To begin with, the learning system has no knowledge about the results of taking a certain action from a certain state, and must take actions more-or-less at random and observe their results in order to acquire knowledge. In a simulated environment, there is no problem with this. However, on a real robot it is possible that this random exploration will place the robot in danger. The canonical example is that a robot does not know that it is bad to try to drive down the stairwell until it has done it, at which point it is too late.

Also, since there are many states and actions to consider, the robot will typically spend all of its early learning time performing random actions and twitching. Only after an extended period of this behavior will it have learned enough to take its first meaningful actions.

One way to safeguard the robot in these early stages of learning is to give it a set of “reflexes” to keep it out of trouble. Whenever the RL system tries to do something dangerous, the reflexes stop it. There are two issues, however, that stem from this. The first is the question of how much knowledge and expertise to give these reflexes. Often for simple learning tasks, once the reflexes are supplied there is little left to learn. In this case, the utility of using a learning approach comes into question.

The second issue is that of coverage of the state space. If the reflexes prevent the robot from entering dangerous parts of the state space, it will never have any experience of these areas. In more sophisticated exploration strategies, the robot will seek out areas of the space about which it has little knowledge. In this case, there will be a conflict; the robot will want to enter the dangerous areas to gain additional knowledge about them, while the reflexes will attempt to prevent it.

These are some of the most important problems facing the use of RL on robots. We now go on to describe a framework for incorporating RL in to robots and discuss how this approach addresses the problems outlined above.

3 The Learning Framework

In this section, we describe a framework for incorporating RL into robots that we are currently working on. The framework has three main components; a provided system, which includes an initial control policy and immediate reward function, a policy learning process and a value learning process. To begin with, the initial control policy has control of the robot. The policy learner passively observes the sensor states and actions taken by the supplied initial policy and attempts to learn the mapping between them. At the same time, the value learner also observes the immediate rewards generated by the supplied policy and attempts to learn a predictive reward function. Once both of these learned mappings are sufficiently accurate, control passes to the learned policy. The system then attempts to improve on this learned policy, based on the the learned predictive reward function.

The provided system consists of an initial control policy and an immediate reward function. The purpose of the initial control policy is to control the robot safely until the learned policy becomes good enough to ensure the robot’s

safety. The implicit assumption is that the task to be learned is such that a reasonable solution to it can be produced in a short period of time. This need not be the optimal solution (or even very close to it), but it should be able to accomplish the goal without placing the robot in any danger.

The immediate reward function is a standard mapping from states to reward for being in that state, as is commonly used in RL. It provides a way of comparing two states quantitatively, and implicitly defines an optimal policy.

3.1 Initial Learning

To begin with, the robot is running the supplied control policy which is assumed to be able to keep the robot out of harm's way, while still making some progress towards the goal. While this is happening, the policy learner observes the sensor values used by the initial policy and the actions generated by it. Using these observations, it constructs the learned policy. After some time, this learned policy will be capable of generating the same behavior as the initial one. At this time, it can be substituted for the initial policy without placing the robot in any danger.

While the policy learner is learning a control strategy, a separate learning process is attempting to learn a predictive reward function. This function will subsequently be used in the second stage of learning, and allows the robot to speculate about the relative merits of different courses of action in any given situation.

Both of the learning processes in this section are supervised learners. That is, they are given both the inputs and the desired outputs and must form the mapping between them. All of the reinforcement learning occurs in the next stage, described below. A key feature of both of these learning algorithms is that they be able to generalize their experiences to instances that they have not seen before. If we make certain continuity assumptions about the forms of the mappings we are trying learn, then it no longer becomes necessary to visit every state-action pair.

3.2 Further Learning

Once the learned policy is generating the same behavior as the original (to within some acceptable accuracy), the system switches control of the robot to the learned policy and begins attempting to improve this mapping. Since the original mapping was not necessarily optimal, it may be improved upon. The criterion for "optimal" and "improvement" is the original supplied reward function. Since this is an immediate reward function, it can only be used to give the reward for an action that was taken. However, the learned predictive reward function can give an estimate of the expected reward for an action that has not yet been taken, and thus reduce the risk of exploratory actions.

On each step in this stage, the learned policy generates an action based on the current sensor state. Based on this action, we generate a number of possible alternative actions. These are then ranked using the predictive reward function, and the one with the greatest expected reward is chosen as the next action for the robot to take. The state-action pair is presented to the policy learner and the action is taken. After the action is complete the immediate reward is calculated and the state-action-reward triplet is presented to the predictive reward learner.

Thus, as time passes, the learned policy changes to more accurately reflect the learned predictive reward function, which itself is changing to reflect the new triplets it sees generated by the learned policy.

4 How the Framework Helps

Perhaps the main contribution of this framework is that the supplied initial control policy bootstraps the learning process and keeps the robot safe during the early stages of learning. It is widely acknowledged that the *tabula rasa* approach to learning, where the learning system starts with no knowledge at all is almost certainly doomed to failure. The initial policy gives the learner an example, however crude, to learn from. By making the first stage of learning supervised (that is, we explicitly provide both the inputs and the outputs), we eliminate the need for dangerous exploration. If we assume that the supplied policy is capable of keeping the robot safe, then we have no need of added "reflexes".

The reward function allows us to take this initially learned policy and refine it further. At the end of the first learning phase, we should be relatively sure of our ability to keep ourselves safe and be able to make reasonable predictions about the consequences of actions, even though we have not taken them from the current state before. Thus, when we wish to generate an exploratory action in order to gain more information about its effects, we can have

some prediction ahead of time. If this prediction suggests that the action will have undesirable consequences (a low, or negative, reward), then it can be avoided. If the generalization ability of our learning algorithms is good, this will let us explore unknown state-action pairs safely without compromising our ability to improve the learned policy.

This framework also allows us to conveniently add bias to the learned policy. Since the initial supervised learning is not concerned with the details of how the outputs are generated, we are free to supply the initial control policy in any form we wish. For example, this can be as a rule-based system, as a large serial program or even as external control, such as using a joystick. It also lets us combine policies from several sources. Imagine, for example, a task that involves following a person down a hallway. We could have modules for obstacle avoidance, corridor traversal and person following running in parallel to control the robot. The initial learned policy will then combine of all of these in a single mapping.

Finally, the framework allows us to specify the final behavior of the robot in more accessible terms. From experience, it seems easier to determine *how* the robot is doing, as opposed to *what* is supposed to be doing. For example it is easier to say “stay away from all obstacles” than it is to code a good obstacle avoidance strategy. In our example, this translates to giving a negative reward whenever the robot is too close to an object.

This allows us to quickly design a solution to the problem and have the RL stage of learning refine it and improve the performance. Again, from our experience, it is relatively easy to write code to do roughly the right thing. The difficulty (and the time-consuming part) is to fine-tune this code. In this framework, we leave this task to the RL system.

5 Issues to be Resolved

At the time of writing we have begun to implement and test the framework described above. We have only very preliminary results but these look somewhat promising. However, our progress to date has highlighted several very important issues that will have to be resolved before the framework becomes of general use.

The most important of these is the choice of supervised learning algorithm. It must have the following properties.

Quick learning Since the amount of input data is likely to be severely limited, we must have an algorithm that can learn and generalize from only a few inputs. It must also learn in real-time. This makes computationally expensive methods, such as backpropagation networks, unappealing.

Able to generalize We cannot hope to visit all of the possible states that the robot might be in or try all of the possible actions. Our learning algorithm should be able to generalize well to unvisited states similar to those it has already seen.

Shifting functions After learning the initial control policy, we attempt to improve on it. This means that the mapping we are trying to learn (from sensor inputs to effector outputs) is changing over time. The learning algorithm should be able to cope with this and be able to determine which previous experiences are still relevant and which are not.

We currently favor instance-based learning methods, in particular a variant of the *k-nearest neighbor* algorithm (*k*-NN). In its most basic form, *k*-NN simply remembers all of the input-output pairs that it has seen. When called upon to provide an output prediction for an unseen input, it finds the *k* most similar instances from its memory and uses these to predict the output for the query point. We are currently experimenting with various type of weighted linear regression for this prediction. By incorporating timestamp information from each training point, we hope to address the problem of shifting output distributions.

One very appealing feature of *k*-NN is its ability to re-evaluate previously seen data when we adopt a new bias. This is most clearly illustrated with an example. Suppose that we weight the inputs according to some scheme, to reflect their relative importance. If, at some point during or after learning, we decide to change these weights, no additional learning is necessary. Since all of the computational expense of *k*-NN is undertaken when a query is processed and also because all experiences are stored in their raw form, we can merely apply the new bias to the old data. Other algorithms, such as backpropagation networks, require extensive retraining in a case like this since they incorporate the bias in their stored representation. The ability to change biases easily is useful and allows us to try out candidates

with very little expense. For example, if we want to try eliminating all but one type of sensor to see if this improves the learned model, then we can do so trivially.

Another problem is that of obtaining suitable coverage of the input space for the learned policy. This is a common problem when dealing with learning from some form of teacher. Typically, the supplied initial control policy will be reasonably competent at completing the task at hand. What this means is that it will keep the robot in that part of the input space that is appropriate for the execution of the task *by that policy*. The robot will, therefore, have never been outside of this envelope of states and consequently will have no idea what to do should it find itself there. This is very bad. Since our robot operates in the real world, events beyond its control might force it into a part of the space in which it has no experience. The best that it can do is then perform randomly, negating any benefits provided by this framework. Note that the supplied initial policy might know what to do in these states should they arise, but since they never arose during training, the learning algorithm was not exposed to them

This leads us to redefine the purpose of the supplied initial control policy. It does not have to be a good solution to the task that we want to learn. Indeed, it can be a very poor solution, often getting into difficulties (and overcoming them) and following indirect paths to the goal. As long as it keeps the robot from harm, this will actually lead to a greater coverage of the input space and more robust behavior during the second learning phase. As an extreme example of this, we attempted to learn a very simple task, where the robot starts at a random distance from a wall and then moves so that it is 2 meters from the wall. Initial supplied policies of a good policy for this task, the “anti-policy” (move in exactly the wrong direction) and a random policy (pick a random direction and move that way until you hit the wall or get to 4 meters away from it) all produced qualitatively the same results during the RL phase of learning.

Another related issue is that of multiple initial policies. A popular way to organize robot control systems is in layers. The lower levels handle basic skills, like obstacle avoidance, while higher layers handle more and more complex behaviors. This can be achieved within our framework by having multiple supervised learning stages. In the first, we learn the most basic skill (perhaps obstacle avoidance). This might involve the robot wandering around randomly avoiding things. As a result, we would cover a large part of the input space. Then we would activate the next higher level of behavior. In the region of the state space where the new policy operates, the previously-learned “default” actions would be superseded by the new ones. This process is repeated for each of the new layers of control, resulting in a single learned policy. This allows the idea of using the framework as a tool for integrating behaviors together. If we have a set of basic behaviors, we can combine them as described above, and then use the supplied reward function to fine-tune them to perform the task that we are interested in.

It is clear that we wish to change from the supervised learning phase to the RL phase at some point after the initial control policy has been learned. However, there are several possible criteria for this change. We do not want to over-fit the data from the initial policy, since this will delay convergence on the final policy. We also do not wish to abandon the supervised stage too early, since this will mean poor models for the RL phase. A related question is whether or not we should keep the initial supplied policy around after we are finished learning from it. If we are called upon to make a prediction in an area of the state space in which we have little experience, it might be better to default to the initial policy than risk a guess from the learned model.

A key feature of the framework is that it can provide predictions on the effects of an unseen action in a given state. The implicit assumption behind this is that the prediction that we ask for is beyond our experience, but only slightly so. This leads to the idea of an area of the state space in which we have had a lot of experience. In this area, we are reasonably sure that our predictions are good. Surrounding this area is a fringe, consisting of states that we have not yet visited, but which are close to those we have. This, our predictions for these states should also be reasonably good, assuming certain continuities in the mappings.

This directly affects our exploration policy. We can no longer select actions at random. Striking a balance between exploration, to gain more knowledge about unknown areas of the state space, and exploitation, taking the best known action, is a difficult problem in RL. It is compounded, by safety issues and also by time taken to generate data, on real robots. We must choose exploratory actions such that their actual outcome is unknown to us but our ability to predict their outcomes is high. By generating a set of such possible actions, we can choose between them in an attempt to maximize the information gain of our exploratory action. This means that we are slowly expanding the area of the state space with which we are familiar until we can be sure that further expansion will be fruitless.

Finally, there is the question of whether or not we should be learning an action model. Currently we learn an action policy, which is a mapping from states to actions, and a predictive reward function, which maps states and actions to

expected reward values. Although it is not necessary to learn an action model, which maps current state and action to next state, doing so greatly facilitates the RL phase of learning. However, learning such a model is often extremely difficult. As stated above, the number of state-action pairs is likely to be huge, requiring us to generalize from the example points that we have seen. However, this imposes continuity assumptions on both the state space and the action space. Often these assumptions are not valid, forcing us to use a more sophisticated (and therefore more expensive) algorithm than we use for the other learners in the framework.

As we noted at the start of this section, we have begun to implement and perform experiments with the framework we describe. We have some initial results from various parts of the system and although these results are very preliminary and qualitative in nature, they seem to support the opinions expressed in this paper.