

Hidden State and Reinforcement Learning with Instance-Based State Identification

R. Andrew McCallum

Department of Computer Science

University of Rochester

Rochester, NY 14627-0226, USA

`mccallum@cs.rochester.edu`

`http://www.cs.rochester.edu/u/mccallum`

Abstract

Real robots with real sensors are not omniscient. When a robot’s next course of action depends on information that is hidden from the sensors because of problems such as occlusion, restricted range, bounded field of view and limited attention, we say the robot suffers from the *hidden state problem*. State identification techniques use history information to uncover hidden state. Some previous approaches to encoding history include: finite state machines [12, 28], recurrent neural networks [25] and genetic programming with indexed memory [49]. A chief disadvantage of all these techniques is their long training time.

This paper presents *instance-based state identification*, a new approach to reinforcement learning with state identification that learns with much fewer training steps. Noting that learning with history and learning in continuous spaces both share the property that they begin without knowing the granularity of the state space, the approach applies instance-based (or “memory-based”) learning to history sequences—instead of recording instances in a continuous geometrical space, we record instances in action-percept-reward *sequence space*.

The first implementation of this approach, called *Nearest Sequence Memory*, learns with an order of magnitude fewer steps than several previous approaches.

I Introduction

This paper is about embedded agents that cannot perceive their whole environment at once. Robotics has seen a surge of interest in reactive robots—agents that use their current sensor values to choose their next action—but as the tasks get more complex, we are finding that some tasks cannot be solved with perception-to-action mappings alone. Sensors are inherently limited, and sometimes the world state information relevant to choosing the next action is hidden by those limitations.

There are many reasons why important features can be hidden from a robot’s perception: sensors have noise, limited range and limited field of view; occlusions hide areas from sensing; limited funds and space prevent equipping the robot with all desired sensors; an exhaustible power supply deters the robot from using all sensors all the time; and the robot has limited computational resources for turning raw sensor data into usable percepts.

In situations where immediate perception does not supply enough information, short-term memory of past perceptions and actions can often provide the information that is needed.

Consider the following examples:

- A space station repair robot sees that it needs a replacement part and thrusts across the station towards a supply bin to get the new part. However, once at the bin, occlusions and limited range prevent the robot from still seeing the distant faulty part. Given only the percepts available from its position at the supply bin, the robot cannot know which part it needs.

With short-term memory, however, the robot could remember which part was faulty, and pick the correct part.

- A hospital delivery robot is making its way from the patient’s room to the blood lab. During its travel, the robot’s limited sensors are not

able to distinguish between many of the different hallway intersections, some in which it should turn right and others in which it should turn left. Given this confusion, the robot cannot reliably deliver the sample.

With short-term memory, however, the robot could distinguish identical-looking intersections by remembering features of its previous locations. For instance, the robot may know that it should turn right at the upcoming intersection because it remembers that it just passed the elevator doors and because the elevator doors are not next to a left-turn intersection.

Unique infrared transceivers installed at the intersections in the hospital might avoid the need for memory, but this installation may not be possible because of another important perceptual limitation listed above: limited funds. The “hardware solution” that uses many transceivers can be considerably more expensive than a “software solution” that uses short-term memory.

- A robot driver uses a vision system to navigate in traffic. In order to overcome limited computational power for extracting usable perceptual information from its raw image data, (*i.e.* the last perceptual limitation in the list above), the driver makes use of an active vision system [3, 40]; it has movable binocular cameras, each containing a fovea. The high resolution in the fovea is necessary for recognizing objects, but if the cameras had such fine-grained resolution over a broad field of view, the vision system would be inundated with far more pixels than its limited computational power could handle.

The driver should pass a car on the highway only when it is close to the preceding car and the blind spot is clear of cars. Unfortunately, addressing the computational limitation has intensified another limitation: limited field of view. The driver cannot see both the car in

front and the blind spot at the same time. Since neither a close-by car ahead nor a clear blind spot alone are enough to tell the robot that it should pull into the passing lane, and since no matter how the robot redirects its active vision system, no immediate perception can include both inputs to the conjunction, the driver cannot decide to pass using immediate perception alone.

With short-term memory, the driver would look at the car ahead, remember that it was close enough for passing, look at the blind-spot for a clear way, and once it had both elements of the relevant conjunction, it could make the decision to begin passing.

A Principles

We propose a “middle path” between purely reactive agents, which suffer from hidden state problems, and traditional planning agents, which require infeasible complete world models. The agent must represent the world state with more than its current perceptions, but the agent does not require a full predictive model of the entire environment.

The work described in this paper addresses the issue of hidden state and short-term memory in conjunction with the following principles:

- Agents should learn their tasks. We would prefer that our robots learn their behaviors rather than have them hard-coded because programming all the details by hand is tedious, we may not know the environment ahead of time, and we want the robot to adapt its behavior as the environment changes.
- Agents should learn in as few trials as possible. Learning experience is expensive in terms of wall clock time, and dangerous in terms of potential damage to the robot and its surroundings. Experience is often

relatively more expensive than storage and computation. Furthermore, the expense of computation and storage will only decrease as advances in computing hardware continue; the cost of experience will not.

- Agents should be able to handle noisy perceptions and actions.
- Strict optimality is not necessary. We should aim for reasonably good behavior given the complexity of the task. In many cases, it's not worth running an algorithm that will take ten times more trials to learn good performance, in order to gain only a one-twentieth increase in the that performance. For example, in a changing environment, this is especially true.

This paper introduces a new family of methods for reinforcement learning with short-term memory. The first implementation of this approach learns with about an order of magnitude fewer steps than several previous approaches. The paper also suggests how future methods in the family could perform even better.

II Background and Related Work

A reinforcement learning¹ agent suffers from hidden state if at any time the agent's state representation is missing information needed to determine the next correct action. More formally, we say a reinforcement learning agent suffers from the *hidden state problem* if the agent's state representation is non-Markovian with respect to actions and utility.

¹This paper assumes that the reader is already familiar with the principles of reinforcement learning as developed in the field of machine learning. The algorithm presented is integrally based on previous work in reinforcement learning, especially *Q*-learning [50]. No review of reinforcement learning basics is included in this paper, but good introductions to reinforcement learning may be found in [50, 5, 52, 30].

The hidden state problem arises as a case of *perceptual aliasing*: the mapping between states of the world and sensations of the agent is not one-to-one [53]. If the agent’s perceptual system produces the same outputs for two world states in which different actions are required, and if the agent’s state representation consists only of its percepts, then the agent will fail to choose correct actions.

Note that even if the agent’s representation consists of more than its percepts—that is, the agent maintains some internal state—it is still possible for the agent to suffer from the hidden state problem. This possibility occurs whenever the agent maintains some amount of internal state, but not enough internal state to disambiguate the aliased world states for the given task.

A Stateless or Fixed-Memory Agents

Some approaches to the hidden state problem do not attempt to disambiguate the aliased states.

The simplest approach is to simply ignore the hidden state problem and apply traditional reinforcement learning methods as if there were no aliased states. Some experience has shown that in certain non-Markovian environments this approach can work [6], however, there are many cases in which ignoring hidden state results in complete failure; see [52, 46, 31] for examples with explanations.

A more careful solution to the hidden state problem is to avoid passing through the perceptually aliased states. This is the approach taken in Whitehead’s *Lion* algorithm [52]. Whenever the agent finds a state that delivers inconsistent reward, it sets that state’s utility so low that the policy will never visit it again. The success of this algorithm depends on a deterministic world and on the existence of a path to the goal that consists of only unaliased states.

Other solutions do not avoid aliased states, but do as best they can given

a non-Markovian state representation while also evading the disasters that would result from ignoring the hidden state altogether [26, 46, 21]. They involve either learning deterministic policies that execute incorrect actions in some aliased states, or learning stochastic policies that, in aliased states, may execute incorrect actions with some probability. These approaches do not depend on a path of unaliased states, but they have other limitations: when faced with many aliased states requiring different actions, the probability of executing an optimal sequence of actions falls precipitously; when faced with potentially harmful results from incorrect actions, deterministically incorrect or probabilistically incorrect action choice may prove too dangerous; and when faced with performance-critical tasks, inefficiency that is proportional to the amount of aliasing may be unacceptable.

B Agents that Learn to Use Memory

The most robust solution to the hidden state problem is to augment the agent’s state representation on-line so as to disambiguate the aliased states. *State identification* techniques uncover the hidden state information—that is, they make the agent’s internal state space Markovian. This transformation from an *imperfect state information model* to a *perfect state information model* has been formalized in the decision and control literature, and involves adding previous percepts and actions to the definition of agent internal state [8]. By augmenting the agent’s percepts with history information—short-term memory of past perceptions, actions and rewards—the agent can distinguish perceptually aliased states, and can then reliably choose correct actions from them.

The example tasks in section 1 illustrated the use of memory to disambiguate hidden state: the robot knows which replacement part to pick because it remembers seeing the faulty part; the robot knows whether to turn right or left at the intersection because it remembers having passed

the elevator door; the robot knows to pull into the passing lane because it remembers having looked both in front and behind. Using memory to disambiguate the aliased situations allows the robot to efficiently perform the task independent of the number of aliased states in that path. It allows the robot to choose correct actions reliably, as opposed to making probabilistic guesses whenever faced with ambiguous sensations.

The problem is how to represent this memory and how the agent should learn what to remember and what to forget. Because the agent begins without knowing how to perform the task at hand, the agent also begins without knowing how much memory will be required.

Predefined, fixed memory representations such as constant-sized perception windows (more formally known as order n Markov models) are often undesirable. When the length of the window is more than needed, they exponentially increase the number of agent internal states for which a policy must be stored and learned; when the length of the memory is less than needed, the agent reverts to the disadvantages of undistinguished hidden state. Even if the agent designer understands the task well enough to know its maximal memory requirements, the agent is at a disadvantage with constant-sized windows because, for most tasks, different amounts of memory are needed at different steps of the task.

Let us emphasize the way in which learning how much to remember is a departure from much other work in machine learning for multi-step tasks. Memory-learning agents must not only learn a mapping from states to actions (which is often difficult enough as it is), but learn the required *state space* as well. When confronted with hidden state, the agent can no longer rely on perception to completely define the agent's internal state space. The internal states of the agent must depend on both the current perception and some variable, learned amount of memory about past perceptions and actions.

C Related Work in Agents that Learn to Use Memory

Several reinforcement learning algorithms are based on techniques that adapt the memory representation on-line. The Perceptual Distinctions Approach [12] and Utile Distinction Memory [31] are both based on splitting states of a finite state machine by doing off-line analysis of statistics gathered over many steps. Recurrent- Q [24] is based on training recurrent neural networks. Indexed Memory [49] uses genetic programming to evolve agents that use load and store instructions on a register bank. A chief disadvantage of all these techniques is that they require a very large number of steps for training.

Before describing these three approaches in more detail, let us consider some other algorithms that are also related.

The environment learning approach implemented in *Toto* [27] also uses memory, and is impressively fast—it learns from single presentations. However, the technique is designed specifically for mobile robots inside rooms and corridors. Unlike the above four algorithms it depends integrally on dead reckoning to disambiguate states, and thus is tied to navigation. In navigational tasks, the global state is directly correlated with geographical position, which is exactly what dead reckoning provides. A meaningful analogy to global state calculating “dead reckoning” for general, non-navigational tasks is not apparent.

ALECSYS [16] is a learning classifier system implemented with a genetic algorithm [20, 9]. In work by Dorigo and Colombetti [13], ALECSYS learns proper sequences by maintaining internal state indicating which subtask the agent should be performing. Transition signals produced by the trainer or environment prompt the agent to switch from one subtask to another. In other work by the same authors [17], ALECSYS is modified to include among its percepts indications of whether a sensor value changed in the last N steps, where N is a parameter of the algorithm set by the robot designer. This is

a variation of the constant-sized perception window approach to memory, with the simplification that the agent remembers only the value since the last change, not all sensor values for the last N steps.

C.1 Finite State Machines

The Perceptual Distinctions Approach (PDA) [12] and Utile Distinction Memory (UDM) [31] both learn *partially observable Markov decision processes* [7] via state-splitting strategies. These probabilistic finite state machines (FSM’s) represent memory in that certain states aren’t likely to be reached except by passing through certain other states previously. So, occupation of the current state can effectively give the agent memory of the perceptions and actions associated with the only available incoming transitions to the current state—the exclusive transitions allow the agent to distinguish states even when the perception associated with those states is identical. The agent can represent arbitrarily long memories by nesting contingent transitions.²

Other research has studied how to calculate an optimal policy given a partially observable Markov decision process [10]. PDA and UDM, on the other hand, both emphasize the learning of the finite state machine itself, then use a simple heuristic for calculating the policy from it. The major difference between PDA and UDM is that UDM only splits states when the

²We should also mention some related work that, while not associated with reinforcement learning, does provide other examples of learning finite state machines.

The map learning work described in [14] learns finite state machines. It is not applicable to hidden state because it requires a unique, unaliased landmark for each world state, (although the perception of the landmarks is allowed to be noisy); also it does not handle stochastically unreliable actions.

The diversity-based inference procedure for finite automata [42, 41] learns an ingenious representation called an *update graph* that can efficiently capture regularities in the structure of the environment. The diversity approach finds a set of canonical tests using repeatable cycles in the environment. The Rivest-Schapire algorithm depends on a completely deterministic environment. Work in a neural network implementation of this approach can accommodate noisy sensations, but has other disadvantages [37].

statistics show that doing so will help predict reward—in other words, UDM only creates as much memory as needed for the task at hand.

The finite state machine approaches have the advantage that there is much established formal study in finite state machines, hidden Markov models and partially observable Markov decision processes. One disadvantage is their limited representational power. For instance, if the agent needs to remember only one bit of information, but must remember that bit for a long time, the entire state space covered during the intervening time must be duplicated. A second disadvantage is that UDM and PDA both require an inordinate number of steps to learn the tasks on which they were demonstrated. They are both based on an Expectation-Maximization (EM) algorithm in which the learner alternates between gathering statistics using the current model, and changing the current model based on the statistics [15]. The agent begins with a random or minimal FSM model, then alternately performs trials consisting of 500 or 1000 steps, and performs state splits or joins based on the statistics from the trial. The large number of steps in each trial is required to make the experience statistically significant. Since, as is always the case with EM algorithms, statistics are based on the current (flawed) model, the requirement of some splits can't be detected until other splits are already made. The result is that only very few splits are made at the end of each trial, many trials are required to complete learning, and at 1000 steps per trial, the learning time suffers dramatically.

C.2 Recurrent Neural Networks

The Recurrent- Q algorithm [24] trains a recurrent neural network containing a fixed number of hidden units. The network holds memory in that the backward looping connections from the hidden units feed into the input layer—the current inputs are partially determined by the networks' own context units, which in turn were determined by previous inputs. Repeated looping of acti-

vation through these recurrent connections can maintain memory over several steps.

An advantage of using recurrent neural networks for memory is that their representational efficiency is greater than finite state machines—they can more efficiently memorize one bit over many steps, for instance. They also have the potential to handle continuous inputs. One disadvantage of Recurrent- Q is that, unlike finite state machines with state splitting rules, the memory capacity is fixed before learning begins since the number of hidden units is fixed. Recurrent- Q also takes very many steps to learn. Neural networks in general are not well known for quick convergence, and the convergence of recurrent neural networks is typically slower than feed forward networks.

C.3 Genetic Algorithms

Genetic Programming with Indexed Memory [49] uses genetic programming [22] to evolve programs that apply load and store instructions to a fixed-sized register bank. These registers, which hold values for later retrieval, obviously supply the agent with memory.

The chief advantage of Genetic Programming with Indexed Memory is its great representational power. In fact, Teller makes a point of emphasizing that the class of potentially evolved programs is Turing complete. Unfortunately, this advantage also causes the algorithm’s worst disadvantage—it pays dearly for this flexibility in terms of learning steps required. Genetic Programming with Indexed Memory can take an extremely large number of steps to learn. For example, in Teller’s demonstration of the algorithm on a grid world box-pushing task, the population contains 800 individuals, the population must be evolved over 100 generations, after each generation all the individuals must be evaluated on 40 test cases that last for 100 steps each—for a total of over 3×10^8 steps.

Our priority is to find a new alternative that learns in fewer steps, with the understanding that we may have to sacrifice flexibility of representation. Agents often cannot afford to take the number of learning steps required for multiple off-line state-splitting tests, or the number of training presentations necessary for a recurrent neural network to converge, or the number of evaluations needed to evolve a population of solutions. Extreme flexibility in memory representation is not always necessary.

III Instance-Based State Identification

This paper advocates a new solution to the hidden state problem we term *instance-based state identification*. The approach was inspired by the successes of instance-based (also called “memory-based”) methods for learning in continuous perception spaces, (*i.e.* [2, 35]). These methods have also had success with physical robots [44, 45].

The term “memory-based” introduces an unfortunate conflict of vocabulary. Here “memory” does not refer to short-term memory of past perceptions and actions used to disambiguate hidden state; it refers to methods that store raw previous experiences directly, such as k -nearest neighbor [18] and locally weighted regression [2]. To avoid confusion, we use the phrase “instance-based,” which we hope carries the same meaning.

The application of instance-based learning to short-term memory for state identification is driven by the important insight that learning in continuous spaces and learning with hidden state have a crucial feature in common: they both begin learning without knowing the final granularity of the agent’s state space. The former learns which regions of continuous input space can be represented uniformly and which areas must be finely divided among many states. The later learns which perceptions can be represented uniformly because they uniquely identify a course of action without the need for memory,

and which perceptions must be divided among many states each with their own detailed history to distinguish them from other perceptually aliased world states. The first approach works with a continuous geometrical input space, the second works with a action-percept-reward “sequence” space, (or “history” space). Large continuous regions correspond to less-specified, small memories; small continuous regions correspond to more-specified, large memories.³

Furthermore, learning in continuous spaces and sequence spaces both have a lot to gain from instance-based methods. In situations where the state space granularity is unknown, it is especially useful to memorize the raw previous experiences. If the agent incorporates experience merely by averaging it into its current, flawed state space granularity, it is bound to attribute experience to the wrong states; experience attributed to the wrong state turns to garbage and is wasted. When faced with an evolving state space, keeping raw previous experience is the path of least commitment, and thus the most cautious about losing information.

Keeping records of raw previous experience does incur certain expenses in terms of storage and computation time, but these expenses can often be justified. For instance, compare the time required to perform a million multiplications with the time it takes the robot to move down the hallway. Often learning trials are more expensive than computation—expensive in many ways: wall clock time, power consumed, danger to the robot and danger to the surroundings, for example. Techniques such as DYNA [48], Experi-

³This analogy between continuous spaces and memory spaces provides a further perspective against implementing memory with constant-sized perception windows: Constant-sized perception windows are undesirable for precisely the same reason constant-sized discretizations of continuous spaces are undesirable. If the discretization is too coarse-grained, the agent will fail; if the discretization is too fine-grained the number of states is exponentially more than needed; different granularity is needed in different regions of state space.

ence Replay [23], Transitional Proximity Q -learning [29] and Asynchronous Dynamic Programming [4] are all examples of efforts to substitute world experience with storage and computation.

IV Nearest Sequence Memory

There are many possible instance-based techniques to choose from, but we wanted to keep the first application simple. With that in mind, this initial algorithm is based on k -nearest neighbor. We call it *Nearest Sequence Memory*, (NSM). It bears emphasizing that this algorithm is the most straightforward, simple, almost naive combination of instance-based methods and history sequences that one could think of; there are still more sophisticated instance-based methods to try. The surprising result is that such a simple technique works as well as it does.

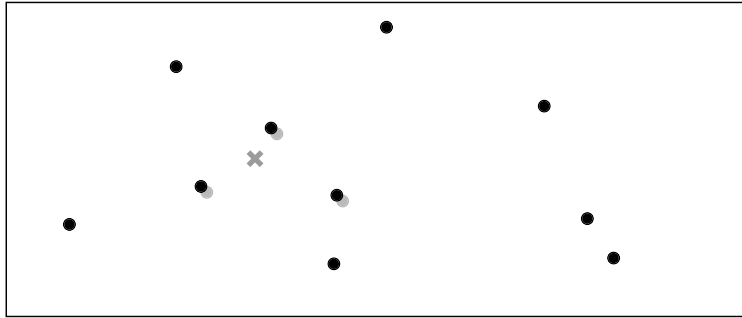
Figure 1 depicts the analogy between k -nearest neighbor in a continuous geometric space and k nearest neighbor in a sequence space.

Any application of k -nearest neighbor consists of three parts: 1) recording each experience, 2) using some distance metric to find neighbors of the current query point, and 3) extracting output values from those neighbors. We apply these three parts to action-percept-reward sequences and reinforcement learning by Q -learning [50] as follows:

1. For each step the agent makes in the world, it records the action, perception and reward by adding a new state to a single, long chain of states. Thus, each state in the chain contains a snapshot of the agent's immediate experience at that moment, and all the agent's experiences are laid out in a time-connected history chain.
2. When the agent is about to choose an action, it finds states considered to be similar by looking in its state chain for states with histories

Learning in a Geometric Space

k -nearest neighbor, $k = 3$



Learning in a Sequence Space

k -nearest neighbor, $k = 3$

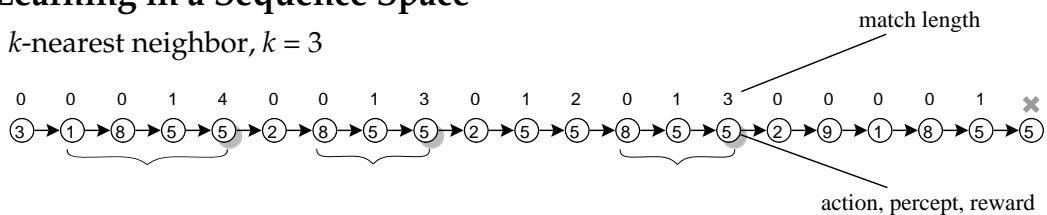


Figure 1: Comparing k -nearest neighbor in a continuous space and in a sequence space. Both have a database of instances, and both have a point which we want to query. In each case, the “query point” from which we measure neighborhood relations is indicated with a gray cross, and the three nearest neighbors are indicated with gray shadows. In a geometric space (shown on top), instances are the black dots, each indicating some vector of values in multi-dimensional space. The neighborhood metric is defined by Euclidean distance. In a sequence space (shown on bottom), instances are circles connected by arrows, each indicating an action/percept/reward snapshot in time, with all the snapshots laid out in time sequence order. (The numbers in the circles indicate action/percept/reward triples.) The neighborhood metric is determined by sequence match length, that is, the number of preceding states that match the states preceding the query point. By using this neighborhood function, the agent finds situations from its past that are most similar to its current situation, where “most similar” not only includes the current percepts but also the percepts leading up to current situation. The agent thus has an internal state space with variable-length short-term memory because when choosing instances from which to extract expected reward values, it tends to prefer instances with better matching history.

similar to the current situation. The longer a state’s string of previous experiences matches the agent’s most recent experiences, the more likely the state represents where the agent is now.

3. Using the states, the agent calculates Q -values by averaging together the expected future reward values associated with the k nearest states for each action. The agent then chooses the action with the highest Q -value. The regular Q -learning update rule is used to update the k states that voted for the chosen action.

Choosing to represent short-term memory as a linear trace is a simple, well-established technique. Order N Markov chains use this representation; the fixed-sized time windows used by Lin and Albus are examples of order N Markov chains [24, 1]. Nearest Sequence Memory also uses a linear trace to represent memory, but it differs from these fixed-sized window techniques because it provides a *variable* memory-length—like k -nearest neighbor, NSM can represent varying resolution in different regions of state space.

V Details of the Algorithm

The interaction between the agent and its environment is described by actions, percepts and rewards. There is a finite set of possible actions, $\mathcal{A} = a^1, a^2, \dots, a^m$, a finite set of possible percepts (observations), $\mathcal{O} = o^1, o^2, \dots, o^n$, and scalar range of possible rewards, $\mathcal{R} = [x, y], x, y \in \mathfrak{R}$. At each time step, t , the agent executes an action, $a_t \in \mathcal{A}$, then as a result receives a new percept, $o_t \in \mathcal{O}$, and a reward, $r_t \in \mathcal{R}$.

We say that the percept “results” from the action in order to be consistent with models of *active* perception in which the percept is a function of both the world state and the previous action [3, 30]. An agent with active perception can execute “perceptual” actions that redirect the agent’s focus of attention.

Note also the atypical grouping of the action, percept and reward in time. In many formulations, the action is given the same time index as the percept and reward that preceded it, *i.e.*, $(o_{t-1}, r_{t-1}, a_{t-1}, o_t, r_t, a_t, o_{t+1}, r_{t+1}, a_{t+1}, \dots)$. The explanation of NSM is made simpler by giving the action the same time index as the percept and reward that *follow* it, *i.e.*, $(o_{t-1}, r_{t-1}, a_t, o_t, r_t, a_{t+1}, o_{t+1}, r_{t+1}, a_{t+2}, \dots)$. In this paper we use the later grouping.

Subscripts on the symbols a , o , and s indicate the time step associated with specific actions, percepts and states. The subscript t refers the the current time step; the subscripts i and j refer to arbitrary time steps. We may add and subtract constants to indicate offsets from the time step indicated by t , i or j .

As with most reinforcement learning, the goal of the agent at each time step is to choose the action that maximizes its expected discounted sum of future reward, called return, and written \mathbf{r}_t .

$$\mathbf{r}_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n r_{t+n} + \dots \quad (1)$$

Just like other instance-based algorithms, Nearest Sequence Memory records each of its raw experiences. The action and environmental state at time t is captured in a “state” data point, written s_t . Recorded with that state is all the available information associated with time t : the action, a_t , the resulting percept, o_t and resulting reward, r_t . Also associated with state s_t is a slot to hold a single expected future discounted reward estimate, denoted $q(s_t)$. This value is associated with action a_t and no other action. We use the notation q to indicate the expected reward value associated with an individual instance; the notation Q refers to the average of the q -values extracted from the k neighboring instances that are closest to the query point.

The Nearest Sequence Memory algorithm consists of the following steps:

1. Find the k nearest neighbor (most similar) states for each possible future action. The state currently at the end of the chain is the “query

point” from which we measure all the distances. The neighborhood metric is defined by “string match” length; it is the number of preceding experience records that match the experience records preceding the “query point” state. This metric can be defined recursively. (Here higher values of $n(s_i, s_j)$ indicate that s_i and s_j are closer neighbors; the notation a_{i-1} indicates the action associated with the state preceding state s_i in the chain.)

$$n(s_i, s_j) = \begin{cases} 1 + n(s_{i-1}, s_{j-1}), & \text{if } (a_{i-1} = a_{j-1}) \wedge (o_{i-1} = o_{j-1}) \wedge (r_{i-1} = r_{j-1}) \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

Considering each of the possible future actions in turn, we find the k nearest neighbors and give them a vote, $v(s_i)$. Ties are resolved by preferring states towards the end of the chain.

$$v(s_i) = \begin{cases} 1, & \text{if } n(s_t, s_i) \text{ is among the } k \text{ } \max_{\forall s_j | a_j = a_i} n(s_t, s_j) \text{'s} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

In the equation above, s_t is the current state, and the “if” expression indicates that a state gets vote 1 only if the neighborhood metric between it and the current state is among the k maximum when compared to the neighborhood similarity values of all the other states that have the same outgoing action.

2. Determine the Q -value for each action by averaging the q -values from the k voting states (states, s , with $v(s) = 1$) from which that action was executed.

$$Q_t(a_i) = \sum_{\forall s_j | a_j = a_i} (v(s_j)/k) q(s_j) \quad (4)$$

3. Select an action by maximum Q -value, or by random exploration. According to an exploration probability, ϵ , either let a_{t+1} be randomly

chosen from \mathcal{A} , or

$$a_{t+1} = \operatorname{argmax}_a Q_t(a) \quad (5)$$

4. Execute the action chosen in step 3, and record the resulting experience. Do this by creating a new “state” representing the current state of the environment, and storing the action-percept-reward triple associated with it:

Increment the time counter: $t \leftarrow t + 1$. Create s_t ; record in it a_t, o_t, r_t .

The agent can limit its storage and computational load by limiting the number of instances it maintains. Instead of allowing the number of instances to grow indefinitely with the number of experience steps, the agent can choose to keep not more than N of them, where N is some reasonably large number like 1000. Once it reaches N instances, it would, with each step, discard the oldest instance and add the new one. The choice of N is balanced by the desire to limit storage and computation, on the one hand, versus the need for the agent to remember experiences from different parts of the environment it visits, on the other. This also provides a way to handle a changing environment.

5. Update the q -values by vote. Perform the dynamic programming step using the standard Q -learning rule to update those states that voted for the chosen action. Note that this actually involves performing steps 1 and 2 to get the next Q -values needed for calculating the utility of the agent’s current state, U_t . (Here β is the learning rate.)

$$U_t = \max_a Q_t(a) \quad (6)$$

$$(\forall s_i | a_i = a_{t-1}) \quad q(s_i) \leftarrow (1 - \beta v(s_i))q(s_i) + \beta v(s_i)(r_i + \gamma U_t) \quad (7)$$

Note the use of r_i instead of r_t in the update of $q(s_i)$. This helps the agent estimate an accurate average reward in a stationary world—using r_t would have biased the average toward recent rewards, where

the strength of the bias would have been dictated by the learning rate, β .

The agent augments the above one-step Q -learning by also passing discounted reward backwards through its recorded history chain for as long as the propagation increases the states' Q -values. In practice, this does not occur often. The process could be described as a kind of “conservative” TD- λ where $\lambda = 1$.

Nearest Sequence Memory performs two kinds of learning. Like all instance-based methods, it learns by populating its input space with raw experience. Unlike most other instance-based methods, it must also perform dynamic programming on values from raw experience. Usually instance-based methods are used as function approximators whose required outputs are the values provided by raw experience—they are supervised learners. In Nearest Sequence Memory, the required output values are Q -values, which are not provided by raw experience, but related to the raw rewards through dynamic programming—Nearest Sequence Memory is a reinforcement learner. These two kinds of learning need not always occur together. For instance, the agent could stop adding states to the chain, but continue to execute steps in the world and do the dynamic programming as implemented by Q -learning. Section 7.2.3 discusses this distinction further.

VI Experimental Results

This section shows experimental results comparing the performance of Nearest Sequence Memory to three other algorithms. Performance is demonstrated using the tasks chosen by the other algorithms' designers. In each case, Nearest Sequence Memory learns the task with roughly an order of magnitude fewer steps. Although Nearest Sequence Memory learns good policies quickly, it does not always learn optimal policies. In section 7.2 we will

discuss why the policies are not always optimal and how Nearest Sequence Memory could be improved.

A Perceptual Distinctions Approach

The Perceptual Distinctions Approach [12] was demonstrated in a space ship docking application with hidden state. The task was made difficult by noisy sensors and unreliable actions. Some of the sensors returned incorrect values 30 percent of the time. Various actions failed 70, 30 or 20 percent of the time, and when they failed, resulted in random states.

The Perceptual Distinctions Approach used learning rate $\beta = 0.1$, temporal discount factor $\gamma = 0.9$, an Expectation-Maximization (EM) cycle of 1000 steps, and an exploration probability $\epsilon = 0.1$. Nearest Sequence Memory used learning rate $\beta = 0.1$, temporal discount factor $\gamma = 0.9$, maximum instance chain length $N = 1000$, exploration probability $\epsilon = 0.1$, and number of neighbors, $k = 8$. A higher k was chosen for this experiment than for the others because of the noisy environment. A graph of the performance during learning appears in figure 2. Performance is indicated using the measure that the agent itself is trying to maximize: sum of future discounted reward (utility). The NSM results have been averaged over five runs. The Perceptual Distinctions Approach takes almost 8000 steps to learn the task. Nearest Sequence Memory learns a good policy in less than 1000 steps, although the policy is not quite optimal.

Not only does NSM learn the task in fewer steps than the Perceptual Distinctions Approach, it also learns the task with less computation. The Baum-Welch procedure [39] is an integral part of the Perceptual Distinctions Approach, and Baum-Welch takes time and space $O(|\mathcal{S}|^2|\mathcal{A}|N)$, where $|\mathcal{S}|$ is the number of states in the model, $|\mathcal{A}|$ is the number of actions and N is the number of steps in the EM cycle. NSM takes time $O(N)$, where N is the maximum instance chain length. For both algorithms N is 1000.

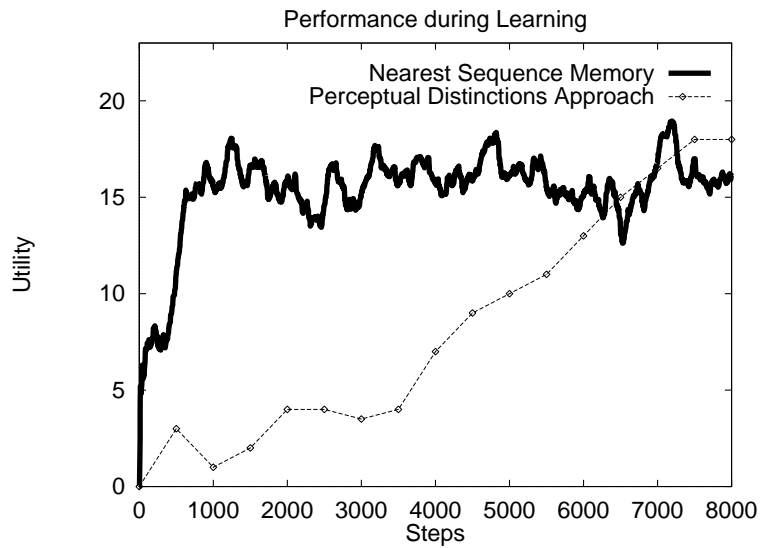


Figure 2: Comparing Perceptual Distinctions Approach and Nearest Sequence Memory on Chrisman’s Docking Task. The spaceship docking task involves noisy sensors and noisy actions. Some sensors returned incorrect values 30% of the time. Some actions failed 70%, 30% or 20% of the time, and when they failed, resulted in random states. The value plotted is the utility (sum of future discounted rewards) of the states the agent’s policy visits.

B Utile Distinction Memory

Utile Distinction Memory [31] was demonstrated on several local perception mazes. Unlike most reinforcement learning maze domains, the agent does not perceive global row and column numbers, but only indications of whether there is a barrier to the immediately adjacent up, down, right or left. Therefore, many maze positions are perceived as identical—including positions from which the agent must execute different actions. Each time the agent reaches the goal, it is reset to a random location in the maze.

Utile Distinction Memory used a learning rate $\beta = 0.9$, temporal discount factor $\gamma = 0.9$, and an Expectation-Maximization cycle of 500 steps. Nearest Sequence Memory used learning rate $\beta = 0.9^4$, temporal discount factor $\gamma = 0.9$, maximum chain length $N = 1000$, exploration probability $e = 0.1$, and number of neighbors $k = 4$. A graph of the number of steps required for learning appears in figure 3. The results are averaged over ten runs. In two of the mazes, Nearest Sequence Memory learns the task in only about 1/20th the time required by Utile Distinction Memory; in the other two, Nearest Sequence Memory learns mazes that Utile Distinction Memory did not solve at all. Although UDM has in principle the capability to learn multi-step memories, in practice they are difficult for UDM to discover [30]. NSM solves these multi-step memory tasks without problem.

Like the Perceptual Distinctions Approach, UDM uses the Baum-Welch procedure. Thus, UDM also has worse computational complexity than NSM.

C Recurrent- Q

Recurrent- Q [24] was demonstrated on a robot two-cup retrieval task. The environment is deterministic, but the task is made difficult by two nested levels of hidden state and by providing no reward until the task is completely

⁴In a deterministic environment the agent can handle a much higher learning rate than in the stochastic environment used above.

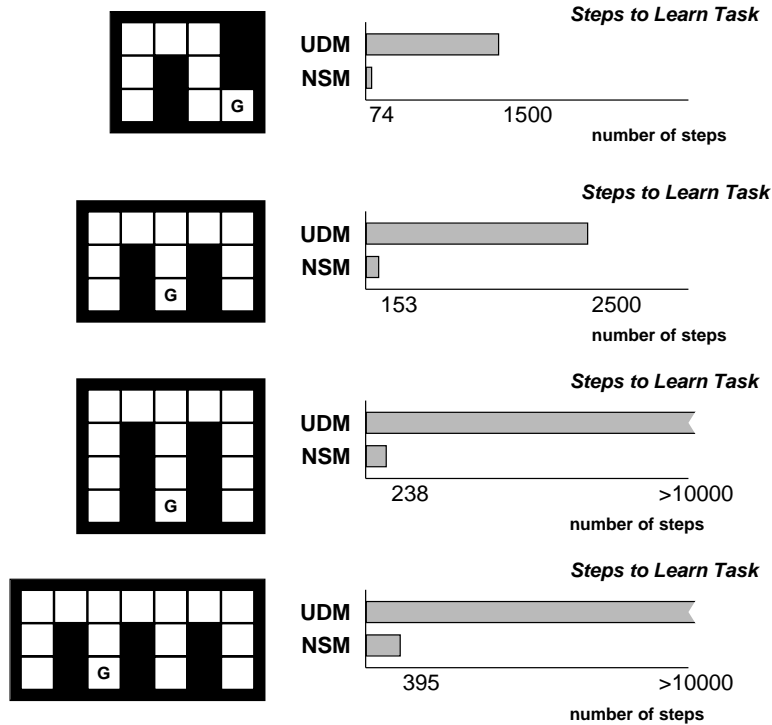


Figure 3: Comparing Utile Distinction Memory and Nearest Sequence Memory on a maze task. The agent does not perceive global row and column numbers, but only indications of whether there is a barrier to the immediately adjacent north, south, east or west. Many positions from which the agent should execute different actions are perceived as identical. The goal square is indicated with a “G”. The bar graphs show the number of training steps required to learn a perfect policy for all starting points. The results are averaged over ten runs.

finished.

Lin actually tried three different neural network algorithms on this task: Window- Q , Recurrent- Q , and Recurrent-Model; Recurrent- Q is the one that learned this task in the fewest trials. Recurrent- Q used temporal discount $\gamma = 0.9$, and a trial time-out of 60 steps. It also used TD- λ ($\lambda = 0.8$) with Experience Replay, which is an experience recording and playback technique developed by Lin for speeding up learning. Nearest Sequence Memory used a trial time-out of 60 steps, learning rate $\beta = 0.9$, temporal discount factor $\gamma = 0.9$, maximum instance chain length $N = 1000$, exploration probability $\epsilon = 0.1$, and number of neighbors $k = 4$. A graph of performance during learning appears in figure 4. Nearest Sequence Memory learns good performance in about 15 trials, Recurrent- Q takes about 100 trials to reach equivalent performance.

A parallel processing implementation of Recurrent- Q could take computation time $O(\text{number of layers in the net})$. A parallel processing implementation of NSM could take computation time $O(1)$. Note also that, like NSM, Lin's Experience Reply mechanism requires the agent to store the chain of raw previous experience.

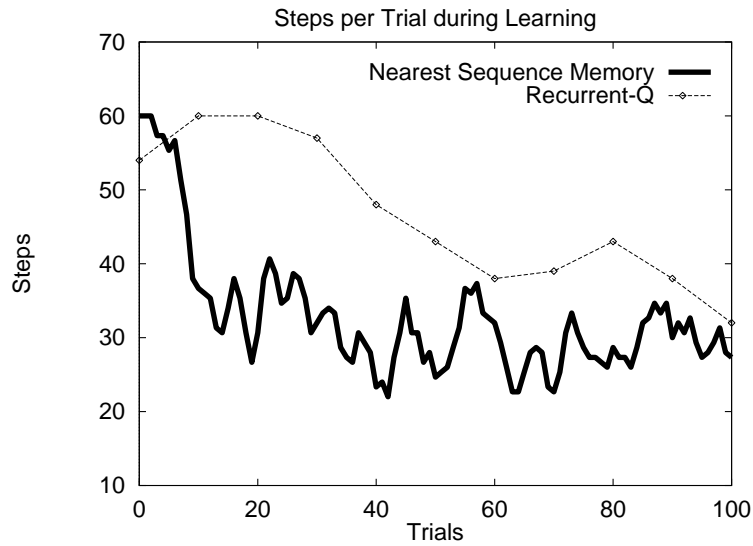


Figure 4: Performance comparison between Recurrent- Q and Nearest Sequence Memory on Lin’s 2-cup task. The robot 2-cup retrieval task is difficult because it has two nested levels of hidden state and because it provides no reward until the task is completely finished. The vertical axis shows the number of steps required to complete the 2-cup task twice, once from each of the two possible starting points. Each of these double task completions makes up a single trial. The horizontal axis measures number of trials since the beginning of learning. The learning curves show the mean results from five runs.

VII Discussion

A Reason for Better Performance

Nearest Sequence Memory offers much improved on-line performance and fewer training steps than its predecessors. Why is the improvement so dramatic? I believe the chief reason lies with the inherent advantage of instance-based state identification, as described in section III, and restated in the following paragraph.

The key idea behind Instance-Based State Identification is the recognition that recording raw experience is particularly advantageous when learning to partition a state space, as is the case when the agent is trying to determine how much history is significant for uncovering hidden state. If, instead of using an instance-based technique, the agent simply averages new experiences into its current, flawed state space model, the experiences will be applied to the wrong states, and cannot be reinterpreted when the agent modifies its state space boundaries. Furthermore, data is always interpreted *in the context of the flawed state space*, always biased in an inappropriate way—not simply recorded, kept uncommitted and open to easy reinterpretation in light of future data.

This scenario is especially clear in the finite state machine approaches: experiences are averaged into the current FSM; after each state-splitting session, the agent re-gathers experience from scratch; the statistics used for learning are derived from the current, flawed FSM; and the raw data used to obtain these statistics is thrown away—not kept for reinterpretation in light of the new state space topology. The same statements apply to the recurrent neural network, except in this case the agent’s internal state space is defined by the weights on connections to the memory-providing hidden units.

B Areas for Improvement

The experimental results in this paper bode well for instance-based state identification. Nearest Sequence Memory is simple—if such a simplistic implementation works as well as it does, more sophisticated approaches may work even better. The following subsections discuss some ideas for improvement.

B.1 Better Distance Metric

The agent should use a more sophisticated neighborhood distance metric than exact string match length. A new metric could account for distances between different percepts instead of considering only exact matches. A new metric could also handle continuous-valued inputs.

There is, however, prior experience showing that neither of these improvements is strictly necessary for work with physical robots. The sequence matching algorithm implemented for *Toto* [27] used a small, finite set of percepts and was completely intolerant to noise. The robot wandered rooms and hallways while avoiding obstacles, learned a map of the environment, and could use the map to reliably travel to arbitrary goal locations. *Toto*'s intolerant sequence matching algorithm overcame its quite noisy sensors by using an intermediate perceptual mechanism that translated the continuous, high-dimensional, noisy sensor values into elements of the useful percept set. Using a little domain knowledge and some averaging over time, the robot's intermediate layer produced noise-less percepts suitable for its intolerant string matching algorithm. A robot using Nearest Sequence Memory could make use of precisely the same mechanism. Furthermore, since NSM is somewhat tolerant of noisy percepts and actions, NSM could even continue to work if the intermediate perceptual mechanism were unreliable.

Unlike finite state machines, Nearest Sequence Memory cannot directly represent cycles in the environment. The agent should use well-chosen or

learned abstract action primitives that encapsulate these loops. Some approach based on variable length, fuzzy sequence matching would also help.

B.2 Structure/Noise Distinction

Nearest Sequence Memory demonstrably solves tasks that involve noisy sensation and action, but it could perhaps handle noise even better if it used some technique for explicitly separating noise from structure.

K -nearest neighbor does not explicitly discriminate between structure and noise [51]. If the current query point has neighbors with wildly varying output values, there is no way to know if the variations are due to noise, (in which case they should all be averaged), or due to fine-grained structure of the underlying function (in which case only the few closest should be averaged). Because Nearest Sequence Memory is built on k -nearest neighbor, it suffers from the same inability to methodically separate history differences that are significant for predicting reward and history differences that are not. I believe this is the single most important reason that Nearest Sequence Memory sometimes did not find optimal policies.

An interesting manifestation of this noise and structure confusion is that occasionally Nearest Sequence Memory exhibits “superstitious behavior” as also displayed by learning animals [19, 38]. In certain circumstances, when the agent first finds success by some round-about route, it tends to use the route repeatedly. When exploratory steps uncover a more efficient path to reward, Nearest Sequence Memory will switch to the better path. For instance, in the beginning of a particular learning run for the Space Ship Docking task, the agent first docked successfully after making two redundant 180 degree turns. The agent then chose this behavior consistently, until, several trials later, random exploration discovered that the two extra turns were unnecessary, and the agent began docking without the extra “dance.” This is over-specification—thinking more detail is relevant than really is. The extra

“dance” is noise, but the agent interprets it as structure. Whether or not the Nearest Sequence Memory exhibits this behavior depends on how densely the state space around the reward is already covered by other instances when it first discovers the reward. This occasional tendency towards extra attention to details does imply that Nearest Sequence Memory could respond well to a teacher, because the agent would learn efficiently from a single presentation, and presumably, a teacher would demonstrate the task without the extra “dances.”

Work in progress addresses the structure/noise issue by combining instance-based state identification with the structure/noise separation method from Utile Distinction Memory [31]. The algorithm, called *Utile Suffix Memory*, uses a tree-structured representation, and is related to work on Prediction Suffix Trees [43], Parti-game [36], G-algorithm [11], and Variable Resolution Dynamic Programming [34]. Preliminary results are very promising. See [32, 33] for more details.

B.3 Separating Instance-Based Learning and Dynamic Programming

As pointed out in the section on details of the algorithm, Nearest Sequence Memory performs two kinds of learning: instance-based population of the state space with raw experience and dynamic programming to calculate expected future reward. Currently these two kinds of learning are intertwined, but I believe the algorithm would perform better with a cleaner separation of the two. The instance-based component should be used to extract the relevant state space, then dynamic programming should be performed only on the extracted states.

The work in progress mentioned in the preceding section makes use of this scheme.

B.4 Further Reducing Storage and Computation

By using the instance-limiting technique discussed in section 5, Nearest Sequence Memory is already computationally more efficient than its competitors. We may want to find ways of reducing the storage and computational requirements even further.

As an alternative to keeping the last N instances, the agent could also simply stop adding states to the chain once it thought it had sufficiently covered the possibly significant sequences. The instance chain would remain unchanged while the agent continued to learn the q -values for those instances.

Once learning is complete, we may want to compile the chain down into a smaller structure. Related work in [47] discusses a technique for compiling example sequences into hidden Markov models. Previous work discusses the application of this technique to building probabilistic finite state machines that include a notion of actions [30].

Acknowledgments

This work has benefited from discussions with many colleagues, including: Dana Ballard, Andrew Moore, Jeff Schneider, and Jonas Karlsson. I am grateful to Dana Ballard, Jonas Karlsson, Virginia de Sa and the IEEE reviewers for making helpful comments on earlier drafts.

This material is based upon work supported by NSF under Grant no. IRI-8903582 and by NIH/PHS under Grant no. 1 R24 RR06853-02.

List of Figures

- 1 Comparing k -nearest neighbor in a continuous space and in a sequence space. Both have a database of instances, and both have a point which we want to query. In each case, the “query point” from which we measure neighborhood relations is indicated with a gray cross, and the three nearest neighbors are indicated with gray shadows. In a geometric space (shown on top), instances are the black dots, each indicating some vector of values in multi-dimensional space. The neighborhood metric is defined by Euclidean distance. In a sequence space (shown on bottom), instances are circles connected by arrows, each indicating an action/percept/reward snapshot in time, with all the snapshots laid out in time sequence order. (The numbers in the circles indicate action/percept/reward triples.) The neighborhood metric is determined by sequence match length, that is, the number of preceding states that match the states preceding the query point. By using this neighborhood function, the agent finds situations from its past that are most similar to its current situation, where “most similar” not only includes the current percepts but also the percepts leading up to current situation. The agent thus has an internal state space with variable-length short-term memory because when choosing instances from which to extract expected reward values, it tends to prefer instances with better matching history. . . . 16

2	Comparing Perceptual Distinctions Approach and Nearest Sequence Memory on Chrisman’s Docking Task. The spaceship docking task involves noisy sensors and noisy actions. Some sensors returned incorrect values 30% of the time. Some actions failed 70%, 30% or 20% of the time, and when they failed, resulted in random states. The value plotted is the utility (sum of future discounted rewards) of the states the agent’s policy visits.	23
3	Comparing Utile Distinction Memory and Nearest Sequence Memory on a maze task. The agent does not perceive global row and column numbers, but only indications of whether there is a barrier to the immediately adjacent north, south, east or west. Many positions from which the agent should execute different actions are perceived as identical. The goal square is indicated with a “G”. The bar graphs show the number of training steps required to learn a perfect policy for all starting points. The results are averaged over ten runs.	25
4	Performance comparison between Recurrent- Q and Nearest Sequence Memory on Lin’s 2-cup task. The robot 2-cup retrieval task is difficult because it has two nested levels of hidden state and because it provides no reward until the task is completely finished. The vertical axis shows the number of steps required to complete the 2-cup task twice, once from each of the two possible starting points. Each of these double task completions makes up a single trial. The horizontal axis measures number of trials since the beginning of learning. The learning curves show the mean results from five runs.	27

References

- [1] J. S. Albus, "Outline for a theory of intelligence," *IEEE Trans. on Systems, Man, and Cybernetics*, vol. 21, no. 3, pp. 473–509, 1991.
- [2] C. G. Atkeson, "Memory-based approaches to approximating continuous functions," in *Nonlinear Modeling and Forecasting* (M. Casdagli and S. Eubank, eds.), pp. 503–521, Addison Wesley, 1992.
- [3] D. H. Ballard and C. M. Brown, "Principles of animate vision," *Computer Vision, Graphics, and Image Processing: Image Understanding*, vol. 56, no. 1, pp. 3–21, 1992.
- [4] A. Barto, S. Bradtke, and S. Singh, "Real-time learning and control using asynchronous dynamic programming," Technical Report 91-57, University of Massachusetts, Amherst, MA, 1991.
- [5] A. B. Barto, R. S. Sutton, and C. J. C. Watkins, "Learning and sequential decision making," in *Learning and Computational Neuroscience* (M. Gabriel and J. W. Moore, eds.), MIT Press, Cambridge, MA, 1990. (Also COINS Tech Report 89-95, Dept. of Computer and Information Sciences, University of Massachusetts, Amherst, MA 01003).
- [6] A. G. Barto, R. S. Sutton, and C. W. Anderson, "Neuron-like elements that can solve difficult learning control problems," *IEEE Trans. on Systems, Man, and Cybernetics*, vol. 13, no. 5, pp. 834–846, 1983.
- [7] D. P. Bertsekas, *Dyanmic Programming and Stochastic Control*. Academic Press, 1976.
- [8] D. P. Bertsekas and S. E. Shreve, *Stochastic Optimal Control*. Academic Press, 1978.

- [9] L. Booker, D. E. Goldberg, and J. H. Holland, “Classifier systems and genetic algorithms,” *Artificial Intelligence*, vol. 40, no. 1-3, pp. 235–282, 1989.
- [10] A. R. Cassandra, L. P. Kaelbling, and M. L. Littman, “Acting optimally in partially observable stochastic domains,” in *Proceedings of the Twelfth National Conference on Artificial Intelligence*, Seattle, WA, 1994.
- [11] D. Chapman and L. P. Kaelbling, “Learning from delayed reinforcement in a complex domain,” in *Twelfth International Joint Conference on Artificial Intelligence*, 1991.
- [12] L. Chrisman, “Reinforcement learning with perceptual aliasing: The perceptual distinctions approach,” in *Tenth National Conference on Artificial Intelligence*, pp. 183–188, 1992.
- [13] M. Colombetti and M. Dorigo, “Training agents to perform sequential behavior,” *Adaptive Behavior*, vol. 2, no. 3, pp. 247–275, 1994.
- [14] T. Dean, D. Angluin, K. Basye, S. Engelson, L. Kaelbling, E. Kokkevis, and O. Maron, “Inferring finite automata with stochastic output functions and an application to map learning,” in *Tenth National Conference on Artificial Intelligence*, pp. 208–214, 1992.
- [15] A. P. Dempster, N. M. Laird, and D. B. Rubin, “Maximum likelihood from incomplete data via the EM algorithm,” *Journal of Royal Statistical Society, Series B*, vol. 39, pp. 1–38, 1977.
- [16] M. Dorigo, “ALECSYS and the autonmouse: Learning to control a real robot by distributed classifier systems,” Technical Report 92-011, Politecnico di Milano, Italy, 1992. (also, *Machine Learning*, in press).

- [17] M. Dorigo and M. Colombetti, “Robot shaping: Developing autonomous agents through learning,” *Artificial Intelligence*, vol. 71, no. 2, pp. 321–370, 1994.
- [18] R. Duda and P. Hart, *Pattern Classification and Scene Analysis*. John Wiley & Sons, New York, NY, 1973.
- [19] E. R. Hilgard and D. G. Marquis, *Conditioning and Learning*. Appleton-Centurn-Crofts, second edition edition, 1961. Edited by Gregory A. Kimble.
- [20] J. H. Holland, *Adaptation in natural and artificial systems*. University of Michigan Press, Ann Arbor, MI, 1975.
- [21] T. Jaakkola, S. P. Singh, and M. I. Jordan, “Reinforcement learning algorithm for partially observable markov decision problems,” in *Advances of Neural Information Processing Systems 7*. Morgan Kaufmann, 1995.
- [22] J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, Cambridge, MA, 1992.
- [23] L. J. Lin, “Self-improvement based on reinforcement learning, planning, and teaching,” in *Proceedings of the Eighth International Workshop on Machine Learning*, 1991.
- [24] L.-J. Lin, *Reinforcement Learning for Robots Using Neural Networks*. PhD dissertation, Carnegie Mellon, School of Computer Science, January 1993.
- [25] L.-J. Lin and T. M. Mitchell, “Reinforcement learning with hidden states,” in *Proceedings of the Second International Conference on Simulation of Adaptive Behavior: From Animals to Animats*, pp. 271–280, 1992.

- [26] M. Littman, “Memoryless policies: Theoretical limitations and practical results,” in *Proceedings of the Third International Conference on Simulation of Adaptive Behavior: From Animals to Animats*, 1994.
- [27] M. J. Mataric, “A distributed model for mobile robot environment-learning and navigation,” Master’s thesis, Massachusetts Institute of Technology, 1990. (Technical Report AI-TR 1228).
- [28] R. A. McCallum, “First results with utile distinction memory for reinforcement learning,” Technical Report 446, University of Rochester Computer Science Dept., 1992.
- [29] R. A. McCallum, “Using transitional proximity for faster reinforcement learning,” in *The Proceedings of the Ninth International Machine Learning Conference*. Morgan Kaufmann Publishers, Inc., 1992.
- [30] R. A. McCallum, “Learning with incomplete selective perception,” Technical Report 453, University of Rochester Computer Science Dept., March 1993. PhD thesis proposal.
- [31] R. A. McCallum, “Overcoming incomplete perception with utile distinction memory,” in *The Proceedings of the Tenth International Machine Learning Conference*. Morgan Kaufmann Publishers, Inc., 1993.
- [32] R. A. McCallum, “Utile suffix memory for reinforcement learning with hidden state,” Technical Report 549, University of Rochester Computer Science Dept., December 1994.
- [33] R. A. McCallum, “Instance-based utile distinctions for reinforcement learning,” in *The Proceedings of the Twelfth International Machine Learning Conference*. Morgan Kaufmann Publishers, Inc., 1995.

- [34] A. W. Moore, “Variable resolution dynamic programming: Efficiently learning action maps in multivariate real-valued state-spaces,” *Proceedings of the Eighth International Workshop on Machine Learning*, pp. 333–337, 1991.
- [35] A. W. Moore, *Efficient Memory-based Learning for Robot Control*. PhD dissertation, University of Cambridge, November 1992.
- [36] A. W. Moore, “The parti-game algorithm for variable resolution reinforcement learning in multidimensional state spaces,” in *NIPS 6*, pp. 711–718. Morgan Kaufmann, 1993.
- [37] M. Mozer and J. Bachrach, “Discovering the structure of a reactive environment by exploration,” *Neural Computation*, vol. 2, pp. 447–457, 1990.
- [38] K. Pryor, *Don’t Shoot the Dog!* Bantam Books, New York, NY, 1984.
- [39] L. R. Rabiner, “A tutorial on hidden Markov models and selected applications in speech recognition,” *Proceedings of the IEEE*, vol. 77, no. 2, February 1989.
- [40] D. A. Reece, *Selective Perception for Robot Driving*. PhD dissertation, School of Computer Science, Carnegie Mellon University, 1992. TR CMU-CS-92-139.
- [41] R. Rivest and R. Schapire, “Diversity-based inference of finite automata,” in *Proceedings of the Twenty Eighth Annual Symposium on Foundations of Computer Science*, pp. 78–86, 1987.
- [42] R. Rivest and R. Schapire, “A new approach to unsupervised learning in deterministic environments,” in *Proceedings of the Fourth International Workshop on Machine Learning* (P. Langley, ed.), pp. 364–375, Irvine, California, 1987.

- [43] D. Ron, Y. Singer, and N. Tishby, "Learning probabilistic automata with variable memory length," in *Proceedings Computational Learning Theory*. Morgan Kaufmann Publishers, Inc., 1994.
- [44] S. Schaal and C. G. Atkeson, "Robot juggling: An implementation of memory-based learning," *Control Systems Magazine*, February 1994.
- [45] J. Schneider, "High dimension action spaces in robot skill learning," in *Twelfth National Conference on Artificial Intelligence (AAAI-94)*, 1994.
- [46] S. P. Singh, T. Jaakkola, and M. I. Jordan, "Model-free reinforcement learning for non-markovian decision problems," in *The Proceedings of the Eleventh International Machine Learning Conference*. Morgan Kaufmann Publishers, 1994.
- [47] A. Stolcke and S. Omohundro, "Hidden Markov model induction by bayesian model merging," in *Advances of Neural Information Processing Systems 4*, November 1992.
- [48] R. S. Sutton, "Integrating architectures for learning, planning, and reacting based on approximating dynamic programming," in *Proceedings of the Seventh International Conference on Machine Learning*, Austin, Texas, 1990, Morgan Kaufmann.
- [49] A. Teller, "The evolution of mental models," in *Advances in Genetic Programming* (K. Kinnear, ed.), ch. , MIT Press, 1994.
- [50] C. Watkins, *Learning from delayed rewards*. PhD dissertation, Cambridge University, 1989.
- [51] D. Wettschereck and T. G. Dietterich, "Locally adaptive nearest neighbor algorithms," in *Advances of Neural Information Processing Systems 6*. Morgan Kaufmann Publishers, Inc., 1994.

- [52] S. D. Whitehead, *Reinforcement Learning for the Adaptive Control of Perception and Action*. PhD dissertation, Department of Computer Science, University of Rochester, 1992.
- [53] S. D. Whitehead and D. H. Ballard, “Learning to perceive and act by trial and error,” *Machine Learning*, vol. 7, no. 1, pp. 45–83, 1991.