
Exploring compact reinforcement-learning representations with linear regression ^{*}

Thomas J. Walsh[†] István Szita[‡]
[†]Dept. of Computer Science
Rutgers University
Piscataway, NJ 08854

Carlos Diuk[†] Michael L. Littman[†]
[‡]Dept. of Computing Science
University of Alberta
Edmonton, AB Canada T6G 2E8

Abstract

This paper presents a new algorithm for on-line linear regression whose efficiency guarantees satisfy the requirements of the KWIK (Knows What It Knows) framework. The algorithm improves on the computational and storage complexity bounds of the current state-of-the-art procedure in this setting. We explore several applications of this algorithm for learning compact reinforcement-learning representations. We show that KWIK linear regression can be used to learn the reward function of a factored MDP and the probabilities of action outcomes in Stochastic STRIPS and Object Oriented MDPs, none of which have been proven to be efficiently learnable in the RL setting before. We also combine KWIK linear regression with other KWIK learners to learn larger portions of these models, including experiments on learning factored MDP transition and reward functions together.

1 Introduction

Linear regression has, for decades, been a powerful tool in the kits of machine-learning researchers. While the field of Reinforcement Learning (RL) [17] has certainly made use of linear regression in approximating value functions [3], using online regression to learn parameters of a model has been limited to environments with linear dynamics (e.g. [7]), and has often been unable

^{*}This Technical Report corrects an earlier version of this work (UAI-09) that erroneously reported an $O(n^2)$ improvement in the KWIK sample complexity of the described linear regression algorithm based on an assumption about the inputs that is not universally valid in the KWIK setting. The reported computational and space improvements do, however, hold. Detailed proofs and a further explanation of the discrepancy are provided in this document.

to make guarantees about the behavior of the resulting learning agent without strict assumptions. One of the great hindrances in applying linear regression to learn models in any RL environment was that the computational and sample efficiency guarantees of online regression learners (such as those that rely on distributional assumptions or do not maintain explicit confidences) did not port to the reinforcement-learning setting, which is not i.i.d. and where realizing what portion of the model needs to be explored is crucial to optimizing reward.

Recently, the introduction of the KWIK (Knows What It Knows) framework [11] has provided a characterization of sufficient conditions for a model-learning algorithm to induce sample-efficient behavior in a reinforcement-learning agent. One of the first algorithms developed for this framework was a KWIK linear regression algorithm [16], which was used to learn the transition function of an MDP with linear dynamics. In this paper, we present an algorithm that improves on both the computational and storage bounds of this previous algorithm and apply it to a stable of learning problems for reinforcement-learning agents that employ “compact” representations. Specifically, we use KWIK linear regression (KWIK-LR) to learn the reward function in a factored MDP and the transition probabilities in domains encoded using Stochastic STRIPS [12] or Object Oriented MDPs (OOMDP) [5]. We note that learning these parameters is not typically associated with linear regression—this paper shows that KWIK-LR can be used to help learn models beyond its standard usage in learning linear dynamics. Because of the KWIK guarantees, agents using this algorithm in these settings are guaranteed to make at most a polynomial (in the parameters of the learned model) number of sub-optimal steps with high probability. We present algorithms and theoretical arguments to this effect, including a general reinforcement-learning algorithm for agents that need to learn the probabilities of action outcomes when these effects may be ambiguous in sample data. Experimental evi-

dence is also presented for benchmark problems in the application areas mentioned above.

Our major contributions are an improved and simplified algorithm for KWIK linear regression, and the sample-efficient algorithms that use KWIK-LR to efficiently learn portions of compact reinforcement-learning representations, none of which have previously been shown to be efficiently learnable. We also discuss extensions, including combining KWIK-LR with other KWIK algorithms to learn more parameters of these compact representations.

2 KWIK Linear Regression

KWIK [11] (Knows What It Knows) is a framework for studying supervised learning algorithms and was designed to unify the analysis of model-based reinforcement-learning algorithms. Formally, a KWIK learner operates over an input space X and an output space Y . At every timestep t , an input $x_t \in X$ is chosen and presented to the learner. If the learner can make an accurate prediction on this input, it can predict \hat{y}_t , otherwise it must admit it does not know by returning \perp (“I don’t know”), allowing it to see the true y_t or a noisy version z_t . An algorithm is said to be KWIK if and only if, with high $(1 - \delta)$ probability, $\|\hat{y}_t - y_t\| < \epsilon$ and the number of \perp s returned over the agent’s lifetime is bounded by a polynomial function over the size of the input problem. It has been shown [10] that in the model-based reinforcement-learning setting, if the underlying model learner is KWIK, then it is possible to build an RL agent around it by driving exploration of the \perp areas using an R-max [4] style manipulation of the value function. Such an agent will, with high probability, take no more than a polynomial (in the parameters of the model being learned) number of suboptimal actions. In this paper, which deals with compact representations where the parameter sizes are far smaller than the enumerated state space, we will call agents that satisfy these conditions, but may require super-polynomial (in the size of the compact parameters) computation for planning, *PAC-compact-MDP*.

One of the first uses of the KWIK framework was in the analysis of an online linear regression algorithm used to learn linear transitions in continuous state MDPs [16]. This algorithm uses the least squares estimate of the weight vector for inputs where the output is known with high certainty. Certainty is measured by two terms representing (1) the number and proximity of previous samples to the current point and (2) the appropriateness of the previous samples for making a least squares estimate. When certainty is low for either measure, the algorithm reports \perp . In this work,

we present a KWIK-LR algorithm that is more sample and computationally efficient than this previous work. First though, we define some notation.

Let $X := \{\vec{x} \in \mathbb{R}^n \mid \|\vec{x}\| \leq 1\}$, and let $f : X \rightarrow \mathbb{R}$ be a linear function with slope $\theta^* \in \mathbb{R}^n$, $\|\theta^*\| \leq M$, i.e. $f(\vec{x}) := \vec{x}^T \theta^*$. Fix a timestep t . For each $i \in \{1, \dots, t\}$, denote the stored samples by \vec{x}_i , their (unknown) expected values by $y_i := \vec{x}_i^T \theta^*$, and their observed values by $z_i := \vec{x}_i^T \theta^* + \eta_i$, where the noise η_i is assumed to form a martingale, i.e., $E(\eta_i \mid \eta_1, \dots, \eta_{i-1}) = 0$, and bounded: $|\eta_i| \leq S$. Define the matrix $D_t := [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_t]^T \in \mathbb{R}^{t \times n}$ and vectors $\vec{y}_t := [y_1; \dots; y_t] \in \mathbb{R}^t$ and $\vec{z}_t := [z_1; \dots; z_t] \in \mathbb{R}^t$, and let I be an $n \times n$ identity matrix.

2.1 A New KWIK Linear Regression Algorithm

Let $X := \{\vec{x} \in \mathbb{R}^n \mid \|\vec{x}\| \leq 1\}$, and let $f : X \rightarrow \mathbb{R}$ be a linear function with slope $\theta^* \in \mathbb{R}^n$, $\|\theta^*\| \leq M$, i.e. $f(\vec{x}) := \vec{x}^T \theta^*$. Fix a timestep t . For each $i \in \{1, \dots, t\}$, denote the stored samples by \vec{x}_i , their (unknown) expected values by $y_i := \vec{x}_i^T \theta^*$, their observed values by $z_i := \vec{x}_i^T \theta^* + \eta_i$, where the noise η_i is assumed to form a martingale, i.e., $E(\eta_i \mid \eta_1, \dots, \eta_{i-1}) = 0$, and bounded: $|\eta_i| \leq S$. Let us define the matrix $D_t := [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_t]^T \in \mathbb{R}^{t \times n}$ and the vectors $\vec{y}_t = [y_1; \dots; y_t] \in \mathbb{R}^t$ and $\vec{z}_t = [z_1; \dots; z_t] \in \mathbb{R}^t$, and let I denote the $n \times n$ identity matrix.

Let $\mathcal{X} := \{\vec{x} \in \mathbb{R}^n \mid \|\vec{x}\| \leq 1\}$, and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be a linear function with slope $\theta^* \in \mathbb{R}^n$, $\|\theta^*\| \leq M$. $f(\vec{x}) := \vec{x}^T \theta^*$. On every timestep $t = 1, 2, \dots$, the agent and the environment interact according to the following protocol:

- the environment poses a query $\vec{x}_t \in \mathcal{X}$,
- the agent either provides an answer $\hat{y}_t \in \mathbb{R}$ or an “I don’t know” answer \perp ,
- If the agent reported \perp ¹ the environment provides the agent a noisy observation $z_t = f(\vec{x}_t) + \eta_t = \vec{x}_t^T \theta^* + \eta_t$. The environment may choose the input \vec{x}_t arbitrarily, and its choice may depend on previous queries and answers $(\vec{x}_1, y_1), \dots, (\vec{x}_{t-1}, y_{t-1})$, however, the observation noises η_t are pairwise independent.

According to [16], an *Admissible KWIK linear regressor* is defined as follows:

¹In the following analysis we make the assumption that the agent actually sees the noisy observation, even if it makes a prediction that is not \perp . While in general this assumption does not fit the KWIK framework, the bounds reported here do not change under this weaker assumption so we use the relaxed (and slightly easier to analyze) protocol in the analysis.

Definition 2.1 We define an admissible algorithm for the KWIK Linear Regression Problem to be one that takes two inputs $0 \leq \epsilon \leq 1$ and $0 \leq \delta < 1$ and, with probability at least $1 - \delta$, satisfies the following conditions: 1. Whenever the algorithm predicts $\hat{y}_t(\vec{x}) \in \mathbb{R}$, we have that $|\hat{y}_t(\vec{x}) - \vec{x}^T \theta^*| \leq \epsilon$. 2. The number of timesteps t for which $\hat{y}_t(\vec{x}_t) = \perp$ is bounded by some function $\zeta(\epsilon, \delta, n)$, polynomial in n , $1/\epsilon$ and $1/\delta$, called the sample complexity of the algorithm.

Let $A_m := \begin{bmatrix} I \\ D \end{bmatrix} \in \mathbb{R}^{(m+n) \times n}$. It is easy to see that the solution of the system $A_t \theta = [(\theta^*)^T; \vec{y}_t^T]^T$ is unique, and equal to θ^* (See Lemma 2). However, the right-hand side of this system is unknown, so we use the approximate system $A_t \theta = [\vec{0}^T; \vec{z}_t^T]^T$, which has a solution $\hat{\theta} = (A_t^T A_t)^{-1} A_t^T [\vec{0}^T; \vec{z}_t^T]^T$. Define $Q_t := (A_t^T A_t)^{-1}$, $\vec{b} := \begin{bmatrix} \theta^* \\ \vec{y} \end{bmatrix} \in \mathbb{R}^{m+n}$, $\vec{c} := \begin{bmatrix} \theta^* \\ \vec{z} \end{bmatrix} \in \mathbb{R}^{m+n}$, $\vec{d} := \begin{bmatrix} \vec{0} \\ \vec{z} \end{bmatrix} \in \mathbb{R}^{m+n}$.

The prediction error on input \vec{x} is

$$\begin{aligned} \hat{y} - y &= \vec{x}^T (\hat{\theta} - \theta^*) \\ &= \vec{x}^T Q_t A_t^T \left(\begin{bmatrix} \vec{0} \\ \vec{z}_t \end{bmatrix} - \begin{bmatrix} \theta^* \\ \vec{y}_t \end{bmatrix} \right) = \vec{x}^T Q_t A_t^T \left(\begin{bmatrix} \vec{0} \\ \vec{\eta}_t \end{bmatrix} - \begin{bmatrix} \theta^* \\ \vec{0} \end{bmatrix} \right). \end{aligned} \quad (1)$$

Algorithm 1 describes our method for KWIK-learning a linear model. The KWIKness of the algorithm is shown by the following theorem:

Theorem 2.2 Suppose that the observation noise is zero-mean and bounded: for all t , $E[\eta_t] = 0$ and $|\eta_t| \leq S$. Let $\delta > 0$ and $\epsilon > 0$. If Algorithm 1 is executed with $\alpha_0 := \min \left\{ c_1 \frac{\epsilon^2}{n}, c_2 \frac{\epsilon^2}{\log \frac{n}{\delta}}, \frac{\epsilon}{2M} \right\}$, with suitable constants c_1 and c_2 , then the number of \perp 's will be $O \left(\max \left\{ \frac{n^3}{\epsilon^4}, \frac{n \log^2 \frac{n}{\delta}}{\epsilon^4} \right\} \right)$, and with probability at least $1 - \delta$, for each sample \vec{x}_t for which a prediction \hat{y}_t is made, $|\hat{y}_t - f(\vec{x}_t)| \leq \epsilon$ holds. ²

In comparison to Strehl & Littman's KWIK online linear regression, the new KWIK-LR requires $\Theta(n^2)$ operations per timestep, in contrast to their $O(tn^2)$ complexity. The new KWIK-LR also requires only $O(n^2)$ space, in contrast to the earlier algorithm, which stored all the samples for which it predicted \perp . Finally, we note that in the case where the inputs X are selected without regard to our agent's behavior (not a fully adversarial situation and therefore not generally applicable in the KWIK setting), it can be shown that the new KWIK-LR algorithm has a KWIK (sample complexity) bound of $O \left(\frac{n}{\epsilon^2} \max \left\{ \frac{1}{\epsilon^2} \log^2 \frac{n}{\delta \epsilon}, M^2 \right\} \right)$,

²The theorem also holds if the noise is Gaussian. There will be a difference only in the constant terms.

Algorithm 1 Learning of linear model

input: α_0
initialize: $t := 0, m := 0, Q := I, \vec{w} := \vec{0}$
repeat
 observe \vec{x}_t
 if $\|Q \vec{x}_t\| < \alpha_0$ **then**
 predict $\hat{y}_t = \vec{x}_t^T Q \vec{w}$ //known state
 else
 predict $\hat{y}_t = \perp$ //unknown state
 observe z_t
 end if
 $Q := Q - \frac{(Q \vec{x}_t)(Q \vec{x}_t)^T}{1 + \vec{x}_t^T Q \vec{x}_t}, \vec{w} := \vec{w} + \vec{x}_t z_t$
 $t := t + 1$
until there are more samples

which is a $\Theta(n^2 / \log^2 n)$ improvement over previously reported KWIK bounds.

3 Proof of Theorem 2.2

3.1 Solution of the regularized system

Lemma 3.1 The least-squares solution of the regularized, noisy, known-optimum system

$$\begin{bmatrix} I \\ D \end{bmatrix} \theta = \begin{bmatrix} \theta^* \\ \vec{y} \end{bmatrix} \quad (2)$$

is unique and equal to θ^* .

Proof. The system has full column-rank, so its least-squares solution can be expressed as

$$\theta = \left(\begin{bmatrix} I \\ D \end{bmatrix}^T \begin{bmatrix} I \\ D \end{bmatrix} \right)^{-1} \begin{bmatrix} I \\ D \end{bmatrix}^T \begin{bmatrix} \theta^* \\ \vec{y} \end{bmatrix} = (D^T D + I)^{-1} (\theta^* + D^T \vec{y}).$$

By reordering, we get that

$$\theta^* = (D^T D + I) \theta - D^T \vec{y} = (D^T D + I) \theta - D^T D \theta^*,$$

or equivalently,

$$(D^T D + I) (\theta^* - \theta) = \vec{0}.$$

The matrix $D^T D + I$ has full rank, so the above equation implies $\theta = \theta^*$. ■

4 Proof of Theorem 2.2

Step 1. We encounter significant uncertainties for a limited number of times only.

Let A_m be our sample matrix before the m -th expansion, and \vec{x}_m be the new sample. Define $Q_m = (A_m^T A_m)^{-1}$. (Note: all A_m have full column ranks, so Q_m is well-defined). We will show that the traces of

the Q_m matrices are positive, monotonously decreasing, and the step sizes are larger than some threshold, consequently, the series has only finitely many elements (which also happens to be bounded by a polynomial of all relevant quantities). The identities of matrix inverses, traces and the statements about their positive definiteness can be found in the Matrix Cookbook [13]. For matrices A and B , we are going to use the notations $A > B$ and $A \geq B$ in the sense that $A - B$ is positive definite or positive semidefinite, respectively. Let $\{\alpha_k\}$ ($k = 0, 1, \dots$) be a monotonously decreasing, 0-limit sequence of constants to be defined later, and let m_k be the number of samples where $\|Q_m \vec{x}_m\| > \alpha_k$ was true.

$Q_1 = I$, so $\text{trace}(Q_1) = n$. For $m \geq 1$, $Q_{m+1} = (Q_m^{-1} + \vec{x}_m \vec{x}_m^T)^{-1} = Q_m - \frac{(Q_m \vec{x}_m)(Q_m \vec{x}_m)^T}{1 + \vec{x}_m^T Q_m \vec{x}_m}$. We need a lower bound on the size of the modification. To this end, note that $A_m^T A_m = (I + D_m^T D_m) \geq I$, so $Q_m = (A_m^T A_m)^{-1} \leq I$. Therefore, $\vec{x}_m^T Q_m \vec{x}_m \leq \vec{x}_m^T \vec{x}_m = \|\vec{x}_m\|^2 \leq 1$. Fix an index $k \geq 0$. The trace decreases significantly in each step when $\|Q_m \vec{x}_m\| > \alpha_k$:

$$\begin{aligned} \text{trace}(Q_{m+1}) &= \text{trace}(Q_m) - \frac{\text{trace}((Q_m \vec{x}_m)(Q_m \vec{x}_m)^T)}{1 + \vec{x}_m^T Q_m \vec{x}_m} \\ &\leq \text{trace}(Q_m) - \frac{1}{2} \text{trace}((Q_m \vec{x}_m)(Q_m \vec{x}_m)^T) \\ &= \text{trace}(Q_m) - \frac{1}{2} \|Q_m \vec{x}_m\|^2 \\ &\leq \text{trace}(Q_m) - \frac{\alpha_k^2}{2} \chi(\|Q_m \vec{x}_m\| > \alpha_k), \end{aligned} \quad (3)$$

where we exploited the fact that for any vector \vec{q} , $\text{trace}(\vec{q}\vec{q}^T) = \text{trace}(\vec{q}^T \vec{q}) = \|\vec{q}\|^2$, and for any matrices Q and R and constant c , $\text{trace}(c \cdot Q) = c \cdot \text{trace}(Q)$ and $\text{trace}(Q + R) = \text{trace}(Q) + \text{trace}(R)$. Applying (3) iteratively, we get that $\text{trace}(Q_{m+1}) \leq \text{trace}(Q_1) - m_k \frac{\alpha_k^2}{2} = n - m_k \frac{\alpha_k^2}{2}$. On the other hand, Q_{m+1} is positive definite, which means that all of its eigenvalues are positive. The trace of a matrix is just the sum of its eigenvalues, so $\text{trace}(Q_{m+1})$ is positive, too. From this, we get that $m_k < 2n/\alpha_k^2$.

Step 2. If we are making a prediction, then it's accurate with high probability.

Suppose we have m samples stored so far. We will drop the subscripts of \vec{x}_m , A_m and Q_m . Suppose that a new sample \vec{x} arrives and the algorithm decides that it is "known", i.e., $\alpha(\vec{x}) \leq \alpha_0$. We have seen that the solution of $A\theta = \vec{b}$ is the optimal weight vector θ^* , that is, $\theta^* = (A^T A)^{-1} A^T \vec{b}$. We will show that for sufficiently large sample size, we will not make large prediction error if we use the solution of $A\theta = \vec{c}$, and ultimately, that is not much different from using the solution of $A\theta = \vec{d}$. Let $\bar{\theta} := (A^T A)^{-1} A^T \vec{c}$. For a new

sample \vec{x} , the error of using $\bar{\theta}$ for prediction instead of θ^* is a weighted sum of the noise:

$$\begin{aligned} \vec{x}^T (\bar{\theta} - \theta^*) &= \vec{x}^T Q A^T (\vec{c} - \vec{b}) \\ &= \vec{x}^T Q [D^T (\vec{z} - \vec{y}) + I^T (\theta^* - \theta^*)] \\ &= \sum_{i=1}^m \vec{x}^T Q \vec{x}_i (z_i - y_i) \\ &= \sum_{i=1}^m (\vec{x}^T Q \vec{x}_i) \eta_i = \vec{x}^T Q A \vec{\eta}. \end{aligned}$$

If $\|Q\vec{x}\| = 0$ then we have $\vec{x}^T \bar{\theta} = \vec{x}^T \theta^*$. Otherwise let k be the index for which $\alpha_{k+1} < \|Q\vec{x}\| \leq \alpha_k$. Note that k is a well-defined index, because $0 < \|Q\vec{x}\| \leq \alpha_0$ according to our assumptions and α_k is monotonously decreasing to 0. Also note that the number of \vec{x} inputs falling in the slot $(\alpha_{k+1}, \alpha_k]$ is at most m_{k+1} . We can bound the squared sum of the weights:

$$\begin{aligned} \sum_{i=1}^m (\vec{x}^T Q \vec{x}_i)^2 &= \sum_{i=1}^m ([\vec{x}^T Q X^T]_i)^2 = \sum_{i=n+1}^{n+m} ([\vec{x}^T Q A^T]_i)^2 \\ &< \sum_{i=1}^{n+m} ([\vec{x}^T Q A^T]_i)^2 = (A Q \vec{x})^T (A Q \vec{x}) \\ &= \vec{x}^T Q (A^T A) Q \vec{x} = \vec{x}^T Q \vec{x} \\ &\leq \|Q\vec{x}\| \|\vec{x}\| \leq \alpha_k. \end{aligned}$$

in the last line we used the Cauchy–Schwarz–Bunyakovski inequality and the fact that $\|Q\vec{x}\| \leq \alpha_k$.

We will use a uniform concentration inequality for martingales. To this end, we need to introduce a couple of notations and definitions. The theorem and the definitions are from chapter 8.2 of Van de Geer's book [19].

Consider a sequence of σ -algebras \mathcal{F}_i and a sequence of parametrized random variables $Z_i(\vec{z})$ having a common parameter \vec{z} , so that $Z_i(\vec{z})$ is \mathcal{F}_i -measurable. Let K be an arbitrary positive real, and $\vec{Z}(\vec{z}) := (Z_1(\vec{z}), \dots, Z_m(\vec{z}))$. Define

$$\begin{aligned} \bar{Z}_m(\vec{z}) &:= \frac{1}{m} \sum_{i=1}^m Z_i(\vec{z}), \\ \bar{A}_m(\vec{z}) &:= \frac{1}{m} \sum_{i=1}^m E(Z_i(\vec{z}) \mid \mathcal{F}_{i-1}), \\ \rho_K(Z_i(\vec{z})) &:= 2K^2 E \left(e^{|Z_i(\vec{z})|/K} - 1 - |Z_i(\vec{z})|/K \mid \mathcal{F}_{i-1} \right), \\ \bar{\rho}_K^2(\vec{Z}(\vec{z})) &:= \frac{1}{m} \sum_{i=1}^m \rho_K(Z_i(\vec{z})) \end{aligned}$$

Definition 4.1 (generalized entropy with bracketing)

Let $0 \leq \delta \leq R$ and fix an integer N . For $1 \leq j \leq N$, let $\vec{L}_j = (L_{1,j}, \dots, L_{m,j})$ and $\vec{U}_j = (U_{1,j}, \dots, U_{m,j})$ be

sequences of random variables so that $L_{i,j}$ and $U_{i,j}$ are \mathcal{F}_i -measurable. Suppose that for all \vec{z} such that $\bar{\rho}_K(\vec{Z}(\vec{z})) \leq R$, there is a non-random index $j = j(\vec{z})$ such that

$$(i). \bar{\rho}_K(\vec{U}_j - \vec{L}_j) \leq \delta^2$$

$$(ii). L_{i,j} \leq Z_i(\vec{z}) \leq U_{i,j} \text{ for all } i.$$

Let $\mathcal{N}_K^b(\delta, R)$ be the smallest non-random value of N for which such a collection exists. Then $\mathcal{H}_K^b(\delta, R) = \log \mathcal{N}_K^b(\delta, R)$ is called the generalized δ -entropy with bracketing.

Theorem 4.2³ Suppose that $\int_0^R \sqrt{\mathcal{H}_K^b(u, R)} du \leq \infty$. For a sufficiently large constant C and for arbitrary constants a, C_0, C_1 that satisfy

$$a \leq C_1 \sqrt{m} R^2 / K,$$

$$a \geq \max \left(C_0 \int_0^R \sqrt{\mathcal{H}_K^b(u, R)} du, C_0 R \right),$$

and

$$C_0^2 \geq C^2(C_1 + 1),$$

we have

$$\Pr \left(\sqrt{m} |\bar{Z}_m - \bar{A}_m| \leq a \vee \bar{\rho}_K(\vec{Z}) \geq R \text{ for all } \vec{z} \right) \geq 1 - C \exp \left(- \frac{a^2}{C^2(C_1 + 1)R^2} \right).$$

To apply the theorem, let $G_k = \{\vec{z} \in \mathbb{R}^n : \|\vec{z}\| \leq \alpha_k\}$, and fix a vector $\vec{z} \in G_k$, and consider the filter defined by the sequence of σ -algebras $\mathcal{F}_i = \sigma(\bar{x}_1, \eta_1, \dots, \bar{x}_i, \eta_i, \bar{x}_{i+1})$. Let $Z_i(\vec{z}) = \vec{z}^T \bar{x}_i \eta_i$. From now on, we will drop \vec{z} from the notation. Note that $|Z_i| \leq S |\vec{z}^T \bar{x}_i| \leq SM \alpha_k$.

Using the fact that $e^x \leq 1 + x + x^2$ for $0 \leq x \leq 1$,

$$\rho_K(Z_i) \leq 2E(Z_i^2 | \mathcal{F}_{i-1}) \leq 2S^2(\vec{z}^T \bar{x}_i)^2$$

for $K \geq SM \alpha_k$, and

$$\bar{\rho}_K^2(\vec{Z}) \leq \frac{2}{m} \sum_{i=1}^m E(Z_i^2 | \mathcal{F}_{i-1}) \leq \frac{2S^2}{m} \|\vec{z}^T A\|^2,$$

³The original theorem had an additional requirement $a \leq 8\sqrt{m}R$, and the integral requirement was

$$a \geq \max \left(C_0 \int_{a/(64\sqrt{m})}^R \sqrt{\mathcal{H}_K^b(u, R)} du, C_0 R \right).$$

However, according to a remark on p.76, if the integral from 0 to R is finite, then we can use it to upper-bound the integral from $a/64\sqrt{m}$ to R , and we can omit the additional condition on a , as its only purpose is to make the lower limit of integration sufficiently small.

The generalized bracketing δ -entropy for G_k is at most $a_1 \cdot n \log \frac{2R}{\delta}$ for some constant $a_1 > 1$, so

$$\begin{aligned} \int_0^R \sqrt{\mathcal{H}_K^b(u, \alpha_k)} du &\leq \sqrt{a_1 \cdot n} \int_0^R \sqrt{\log \frac{2R}{u}} du \\ &\leq \sqrt{a_1 \cdot n} \int_0^{2R} \sqrt{\log \frac{2R}{u}} du \\ &= \sqrt{a_1 \pi \cdot n} R = a_2 \sqrt{n} R, \end{aligned}$$

where we implicitly introduced a new constant $a_2 > 1$. Let us fix $R = \sqrt{\frac{\alpha_k M}{m}}$, $K = SM \alpha_k$, $a = \frac{\epsilon_1}{\sqrt{m}}$, $C_0 = \sqrt{S+1}C$, and $C_1 = S$. With these assignments, the inequality between C, C_1 and C_0 is automatically satisfied, and the other two inequalities reduce to

$$\sqrt{\alpha_k n M (S+1)} a_2 C \leq \epsilon_1 \leq 1.$$

An equivalent form of the left-hand-side inequality is

$$\alpha_k \leq \frac{\epsilon_1^2}{n M (S+1) a_2^2 C^2}.$$

Furthermore, we would like the failure probability to be less than $\delta_1/2^k$. This holds if

$$\alpha_k \leq \frac{\epsilon_1^2}{C^2 (S+1) M \log \frac{2^k}{C \delta_1}}.$$

By these assignments, the statement of the theorem becomes

$$\Pr \left(|\vec{z}^T A \vec{\eta}| \leq \epsilon_1 \vee \|A \vec{z}\|^2 \geq \alpha_k M \text{ for all } \vec{z} \right) \geq 1 - \delta_k.$$

Specifically, for $\vec{z} = Q\vec{x}$, we have $\|AQ\vec{x}\|^2 \leq \alpha_k M$, so

$$\Pr \left(|\vec{x}^T Q A \vec{\eta}| \leq \epsilon_1 \right) \geq 1 - \delta_k.$$

We now show that when a sample is “known”, it does not really matter, what target values are given for the first n coordinates. Specifically, we can use the $\vec{0}$ vector instead of θ^* , without committing a large error. Let $\hat{\theta} := (A^T A)^{-1} A^T \vec{d}$. The error made by using $\hat{\theta}$ instead of θ^* is

$$\begin{aligned} \vec{x}^T (\hat{\theta} - \theta^*) &= \vec{x}^T Q A^T (\vec{d} - \vec{c}) \\ &= \vec{x}^T Q \begin{bmatrix} I \\ X \end{bmatrix}^T \begin{bmatrix} \theta^* \\ \vec{0} \end{bmatrix} = \vec{x}^T Q \theta^*. \end{aligned}$$

which can be bounded by

$$\left| \vec{x}^T (\hat{\theta} - \theta^*) \right| = \left| \vec{x}^T Q \theta^* \right| \leq \|Q\vec{x}\| \|\theta^*\| \leq \alpha_0 M.$$

This will be smaller than ϵ_1 if $\alpha_0 < \epsilon_1/M$.

Let us substitute $\epsilon_1 = \epsilon/2$. Then we have that if

$$\alpha_k := \min \left\{ \frac{\epsilon^2}{4nM(S+1)a_2^2 C^2}, \frac{\epsilon^2}{4C^2(S+1)M \log \frac{2^k}{C \delta_1}}, \frac{\epsilon}{2M} \right\},$$

then $|\hat{y} - y| < \epsilon$ with probability at least $1 - \delta_1/2^k$:

$$\begin{aligned} |\hat{y} - y| &= \left| \bar{x}(\hat{\theta} - \theta^*) \right| \\ &\leq \left| \bar{x}(\hat{\theta} - \bar{\theta}) \right| + \left| \bar{x}(\bar{\theta} - \theta^*) \right| \leq \frac{\epsilon}{2} + \frac{\epsilon}{2}. \end{aligned}$$

With this setting of α_0 , the number of “don’t know”s is at most

$$\begin{aligned} \bar{m} = \frac{2n}{\alpha_0^2} &:= \max \left\{ \frac{32n^3 M^2 (S+1)^2 a_2^4 C^4}{\epsilon^4}, \right. \\ &\left. \frac{32n C^4 (S+1)^2 M^2 \log^2 \frac{2}{C\delta_1}}{\epsilon^4}, \frac{8nM^2}{\epsilon^2} \right\}. \end{aligned}$$

Furthermore, the sequence $\{\alpha_k\}$ is monotonously decreasing, as required, and the total probability of an error is at most

$$\begin{aligned} &\sum_{k=0}^{\infty} m_k \frac{\delta_1}{2^k} \leq \sum_{k=0}^{\infty} \frac{2n}{\alpha_k^2} \frac{\delta_1}{2^k} \\ &\leq \delta_1 \max \left\{ \frac{32n^3 M^2 (S+1)^2 a_2^4 C^4}{\epsilon^4} \sum_{k=0}^{\infty} \frac{1}{2^k}, \right. \\ &\quad \frac{32n C^4 (S+1)^2 M^2}{\epsilon^4} \sum_{k=0}^{\infty} \frac{\log^2(2^{k+1}/C\delta_1)}{2^k}, \\ &\quad \left. \frac{8nM^2}{\epsilon^2} \sum_{k=0}^{\infty} \frac{1}{2^k} \right\} \\ &= \delta_1 \max \left\{ \frac{64n^3 M^2 (S+1)^2 a_2^4 C^4}{\epsilon^4}, \right. \\ &\quad \frac{32n C^4 (S+1)^2 M^2}{\epsilon^4} (12 \log^2 2 + \\ &\quad 4 \log \frac{2}{C\delta_1} \log 2 + 2 \log^2 \frac{2}{C\delta_1}), \\ &\quad \left. \frac{16nM^2}{\epsilon^2} \right\} \\ &= \Theta \left(\max \left\{ \frac{n^3 \delta_1}{\epsilon^4}, \frac{n \delta_1 \log^2(1/\delta_1)}{\epsilon^4}, \frac{n \delta_1 M^2}{\epsilon^2} \right\} \right), \end{aligned}$$

where we used the identities $\sum_{k=0}^{\infty} 1/2^k = 2$, $\sum_{k=0}^{\infty} k/2^k = 2$ and $\sum_{k=0}^{\infty} k^2/2^k = 6$. Let us set δ_1 so that the above probability is less than δ . For this, the assignment

$$\delta_1 = O \left(\min \left\{ \frac{\epsilon^4 \delta}{n^3}, \frac{\epsilon^4 \delta^2}{n}, \frac{\epsilon^2 \delta}{nM^2} \right\} \right)$$

is sufficient. \blacksquare

5 Application 1: Learning Rewards in a Factored-State MDP

A Markov decision process (MDP) [17] is characterized by a quintuple $(\mathbf{X}, A, R, P, \gamma)$, where \mathbf{X} is a finite set of states; A is a finite set of actions;

$R : \mathbf{X} \times A \rightarrow \mathbb{R}$ is the reward function of the agent; $P : \mathbf{X} \times A \times \mathbf{X} \rightarrow [0, 1]$ is the transition function; and finally, $\gamma \in [0, 1]$ is the discount rate on future rewards. A factored-state Markov decision process (fMDP) is a structured MDP, where \mathbf{X} is the Cartesian product of m smaller components: $\mathbf{X} = X_1 \times X_2 \times \dots \times X_m$. A function f is a *local-scope* function if it is defined over a subspace $\mathbf{X}[Z]$ of the state space, where Z is a (presumably small) index set. We make the standard assumption [9] that for each i there exist sets Γ_i of size $O(1)$ such that $\bar{x}_{t+1}[i]$ depends only on $\bar{x}_t[\Gamma_i]$ and a_t . The transition probabilities are then $P(\bar{y} | \bar{x}, a) = \prod_{i=1}^m P_i(\bar{y}[i] | \bar{x}[\Gamma_i], a)$ for each $\bar{x}, \bar{y} \in \mathbf{X}$, $a \in A$, so each factor is determined by a local-scope function. We make the analogous assumption that the reward function is the sum of J local-scope functions with scopes Z_j : $R(\bar{x}, a) = \sum_{j=1}^J R_j(\bar{x}[Z_j], a)$. An fMDP is fully characterized by the tuple $\mathcal{M} = (\{X_i\}_{i=1}^m; A; \{P_i\}_{i=1}^m; \{R_j\}_{j=1}^J; \bar{x}_s; \gamma; \{\Gamma_i\}_{i=1}^m; \{Z_j\}_{j=1}^J)$. Algorithms exist for learning transition probabilities [9] and dependency structure [15], but until now, no algorithm existed for learning the reward functions.

For $J > 1$, we can only observe the *sums* of unknown quantities, not the quantities themselves. Doing so requires the solution of an online linear regression problem with a suitable encoding of the reward model. Let N_r be the total number of parameters describing the reward model ($N_r \leq J|A|n^{n_f}$), and consider the indicator function $\chi(\bar{x}, a) \in \mathbb{R}^{N_r}$, which is 1 on indices corresponding to $R_j(\bar{x}[Z_j], a)$ and 0 elsewhere.

Our solution, Algorithm 2, is a modification of Algorithm 1. We initialize all unknown rewards to some constant R_0 (analogous to the common maximum reward parameter R_{\max} [4]). If R_0 is sufficiently high, the algorithm outputs optimistic reward estimates for unknown states (instead of \perp), and otherwise gives near-accurate predictions with high probability. This property follows from standard arguments [18]: for unknown states, the noise term of the prediction error can be bounded by Azuma’s inequality, and R_0 can be set high enough so that the second term is positive and dominates. This form of optimistic initialization has proven consistently better than R-max in flat MDPs [18]. For known states, the KWIK guarantees suffice to ensure near-optimal behavior. Because it combines a KWIK model-learner with R-max style exploration, this algorithm is PAC-compact-MDP—the first efficient algorithm for this task. In Section 7.1 we combine this algorithm with another KWIK learner that learns the transition dependency structure and probabilities [11] to learn the full fMDP model.

Algorithm 2 Reward learning in fMDPs with optimistic initialization

input: R_0
 $\mathcal{M}^{-r} = (\{X_i\}_1^m; A; \{P_i\}_1^m; \vec{x}_s; \gamma; \{\Gamma_i\}_1^m; \{Z_j\}_1^J)$
 $t := 0, \vec{x}_0 := \vec{x}_s, m := 0, Q := I, \vec{w} := R_0 \vec{1}$
repeat
 $\hat{r} := Q\vec{w}$
 $a_t :=$ greedy action in fMDP($\mathcal{M}^{-r}, \hat{r}$), state \vec{x}_t
 execute a_t , observe reward r_t , next state \vec{x}_{t+1}
 $\vec{u} := \chi(\vec{x}_t, a_t)$
 $Q := Q - \frac{(Q\vec{u})(Q\vec{u})^T}{1 + \vec{u}^T Q \vec{u}}, \vec{w} := \vec{w} + \vec{u}^T r_t$
 $t := t + 1$
until there are more samples

5.1 Experiments

We carried out reward-learning experiments on the Stocks domain [15], with 3 sectors and 2 stocks per sector. Rewards were uniformly random in the interval $[0.5, 1.5]$ for owning a rising stock, and random in $[-1.5, -0.5]$ for owning a decreasing stock. We compared six algorithms: (1) Algorithm 2; (2) Algorithm 1 modified to output R_{\max} in unknown states; (3) the previous state-of-the-art KWIK-LR algorithm [16] modified to output R_{\max} in unknown states; (4) a flat tabular reward learner; and to demonstrate the need for efficient exploration, (5) linear regression without exploration and (6) linear regression with epsilon-greedy exploration.

Each algorithm was run 20 times for 250 steps, updating the model every 5 steps. For the Stocks(3,2) domain, $R_{\max} = 6$. We used $R_0 = 10$ for Algorithm 2, $\alpha_0 = 1.0$ for the second approach and $\alpha_1 = \alpha_2 = 1.0$ for the third one. The values of the learned policies are shown in Figure 1. All curves except (1) and (2) differ significantly on a 95% confidence level. Notice (3) and (5) take longer to learn the model, (4) takes far longer, and (6) fails to explore and find the correct model.

6 Learning Transition Probabilities

We consider another novel application of KWIK linear regression—learning action-effect probabilities in environments where these effects may be ambiguous. Specifically, we consider environments where actions are encoded as stochastic *action schemas* (e.g. travel(X, Y) rather than travel(Paris, Rome)) and the effects of these actions are stochastic. For instance, the action travel(X, Y) may result in the effect $at(Y)$ with probability .9 and the effect $at(X)$ with probability .1. More formally, every action $a \in A$ is of the form $a = [(\omega_0, p_0) \cdots (\omega_n, p_n)]$ where each $\omega_i \in \Omega^a$ is

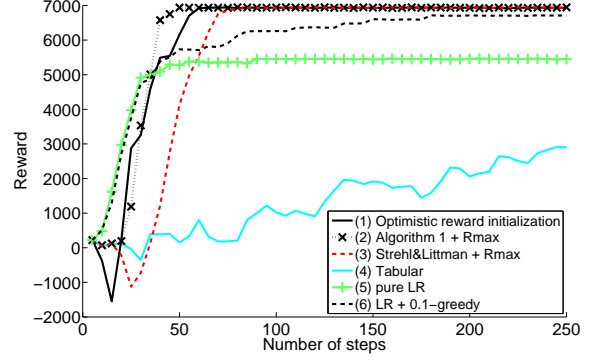


Figure 1: The value of the learned policy as a function of time for different reward-learning algorithms.

Algorithm 3 Transition Probability Learner

1: input: S, A (action schemas sans p_i s), R, α
2: $\forall a \in A$, instantiate a learner (Algorithm 1) $\mathcal{L}^a(\alpha)$
3: for each current state s_t **do**
4: **for** each $e \in E(s, a)$, $s \in S$ and $a \in A$ **do**
5: Construct \vec{x} where $x_j = 1$ if $j \in e$, else 0.
6: $\hat{P}(e) =$ Prediction of $\mathcal{L}^a(\vec{x})$ // can be \perp
7: end for
8: Perform modified value iteration on $\{S, A, E, \hat{P}, R\}$, to get greedy policy π
9: Perform action $a_t = \pi(s_t)$, observe $e_t \in E(s_t, a_t)$
10: for equivalence classes $e \in E(s_t, a_t)$ **do**
11: Construct \vec{x} where $x_j = 1$ if $j \in e$, else 0.
12: Update \mathcal{L}^{a_t} with \vec{x} and $y = 1$ if $e = e_t$, else $y = 0$
13: end for
14: end for

a possible effect. When the action is taken, one of these effects occurs according to the probability distribution induced by the p_i s. The schemas may also contain conditional effect distributions, the nature of which is determined by the specific language used (as discussed in the following subsections). This form of generalization has been used to encode many different types of environments in RL, including stochastic STRIPS [12], Object Oriented MDPs (OOMDPs) [5], and typed dynamics [14]. Learning these probabilities is non-trivial because for a given state action pair (s, a) , the effects are partitioned into equivalence classes $E(s, a) = \{\{\omega_i, \omega_j, \omega_k\}, \{\omega_l, \omega_m\}, \dots\}$ where each $e \in E(s, a)$ contains effects that are identical given state s . For instance, if an action has two possible effects $Holding(X)$ and $Holding(Y)$, but is taken in a state where X and Y are both already held, we cannot tell which actually occurred. Notice that the probability of any equivalence class is equal to the sum of the probabilities of the effects it contains, hence the link

to linear regression.

The standard “counting” method for learning probabilities cannot be used in such a setting because it is unclear which effect’s count (*Holding(X)*s or *Holding(Y)*s) we should increment. However, KWIK-LR can be used to learn the probabilities. While standard linear regression could also be used if transition data was available in batch form, an RL agent using KWIK-LR can learn the probabilities online and maintain the PAC-compact-MDP guarantee needed for effective exploration. Algorithm 3 presents a PAC-compact-MDP algorithm for an agent that is given full action-operator specifications, except for the probabilities. That is, since our focus is on learning the probabilities, we assume that each ω_i is known, but not the p_i s. We also assume the preconditions and reward structure of the problem are known. We discuss methods for relaxing these assumptions in Section 7. This algorithm can be used with several representations including those mentioned above.

Intuitively, Algorithm 3 computes an optimistically optimal action a_t using a planner (Line 7, detailed below) and then gathers experience indicating which of the equivalence classes $e_t \in E(s_t, a_t)$ actually occurred. For instance, it may see that *Holding(X)* occurred, instead of *Holding(Y)*. This equivalence class contains one or more effects ($\omega_i, \omega_j \dots \in \Omega^a$), and an indicator vector \vec{x} is created where $\vec{x}_i = 1$ if $\omega_i \in e_t$ (Line 10). For instance, if the agent was not holding anything, and then *Holding(X)* occurred and not *Holding(Y)*, the vector would be $[1; 0]$, but in a state where both were already held, \vec{x} would always come out $[1; 1]$. Note that each equivalence class in $E(s, a)$ induces a unique (and disjoint) \vec{x} . Each of these is used to update the KWIK-LR learner (Lines 9-11), with an output (y) of 1 associated with the \vec{x} that actually happened (which may have 1s in multiple places, as in the ambiguous *Holding* case). Given any possible action and equivalence class in the state/action space, the KWIK-LR agent can now be queried to determine the probability of the equivalent transition (Lines 5-6), though it may return \perp , identifying transitions from equivalence partitions that have unlearned probability.

Planning in Algorithm 3 is done by constructing a transition model in the grounded state space (Lines 4-6). KWIK-LR determines for each $E(s, a)$, the probability of the possible next states (one for each $e \in E(s, a)$). A modified version of value iteration [17] (Line 7) is then used to plan the optimal next action. The modification is that at every iteration, for every state-action pair that has effects with known probabilities $K = \{\omega_i, \omega_j \dots\}$ and unknown probabilities $U = \{\omega_k, \omega_l \dots\}$, the effect in U that leads to the highest value next state is considered to have probability

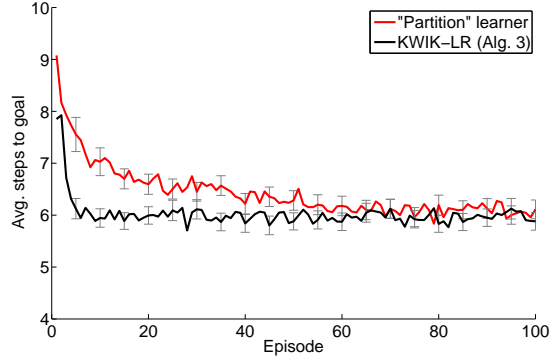


Figure 2: Average steps to the goal in Stochastic Paint/Polish World over 1000 runs.

| |
|---|
| <p><i>paint(X)</i>: reward = -1 PRE: none ADD: <i>Painted(X)</i> DEL: none 0.6 ADD: <i>Painted(X)</i>, <i>Scratched(X)</i> DEL: none 0.3 ADD: none DEL: none 0.1</p> <p><i>polish(X)</i>: reward = -1 PRE: none ADD: none DEL: <i>Painted(X)</i> 0.2 ADD: none DEL: <i>Scratched(X)</i> 0.2 ADD: <i>Polished(X)</i> DEL: <i>Painted(X)</i>, <i>Scratched(X)</i> 0.3 ADD: <i>Polished(X)</i> DEL: <i>Painted(X)</i> 0.2 ADD: none DEL: none 0.1</p> <p><i>shortcut(X)</i>: reward = -1 PRE: none ADD: <i>Painted(X)</i>, <i>Polished(X)</i> DEL: none 0.05 ADD: none DEL: none 0.95</p> <p><i>done(X)</i>: reward = 10 PRE: <i>Painted(X)</i>, <i>Polished(X)</i>, <i>Scratched(X)</i> ADD: <i>Finished(X)</i> DEL: none 1.0</p> |
|---|

Table 1: Stochastic Paint/Polish world

$1 - \sum_{\omega_i \in K} P(\omega_i)$. This change is designed to force value iteration to be optimistic—it considers the most rosie of all models consistent with what has been learned. We note that planning in the flat state space can require exponential computation time, but this is often unavoidable, and since the model is still *learned* in a sample efficient manner, it satisfies the conditions for PAC-compact-MDP.

6.1 Application 2: Stochastic STRIPS

STRIPS domains [8] are made up of a set of objects \mathcal{O} , a set of predicates P , and a set of actions A . The actions have conjunctive (over P) preconditions and effects specified by **ADD** and **DELETE (DEL)** lists, which specify what predicates are added and deleted from the world state when the action occurs. Stochastic STRIPS operators generalize this representation

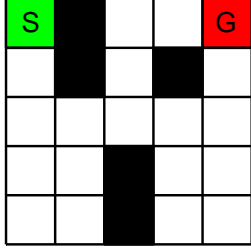


Figure 3: Maze domain.

by considering multiple possible action effects specified by $\langle \text{ADD}, \text{DEL}, \text{PROB} \rangle$ tuples as in Table 1. Notice this representation is an instance of the general action schemas defined above. While others have looked at learning similar operators [12], their work attempted to heuristically learn the full operators (including structure), and could not give any guarantees (as we do) on the behavior of their algorithm, nor did they identify an efficient algorithm for learning the probabilities, as we have with Algorithm 3. To make the planning step well defined, we consider Stochastic STRIPS with *rewards* [20].

Table 1 illustrates a variant of the familiar Paint/Polish domain in the Stochastic STRIPS setting. There are several ambiguous effects. For instance, executing *paint* on a scratched but not painted object, and observing it is now scratched and painted, one cannot tell which of the first two effects occurred. We used this domain with a single object to empirically test Algorithm 3 against a version of the algorithm that does not use KWIK-LR, instead attempting to learn every possible equivalence class partition distribution separately (*Partition*). Because the focus here is on learning the transition probabilities ($p_{i,s}$), both learners were given the preconditions and effects (ω_i) of each action. We discuss relaxing these assumptions later. Figure 2 shows the results (with the known/unknown thresholds empirically tuned) averaged over 1000 runs with randomized initial states for each episode (both learners receive the same initial states). Algorithm 3 learns much faster because it effectively shares information between partitions.

6.2 Application 3: Stochastic Object Oriented RL

Object-oriented MDPs [5] consist of a set of objects \mathcal{O} , a set of actions A that take elements of \mathcal{O} as parameters, and (in the original deterministic description) a set of condition-effect pairs $\langle c, \omega \rangle$ associated with each action. Objects have attributes, and the set of all objects and their attribute values at a given time t constitute the state s_t . When an action a is executed from s_t , the environment checks which condition

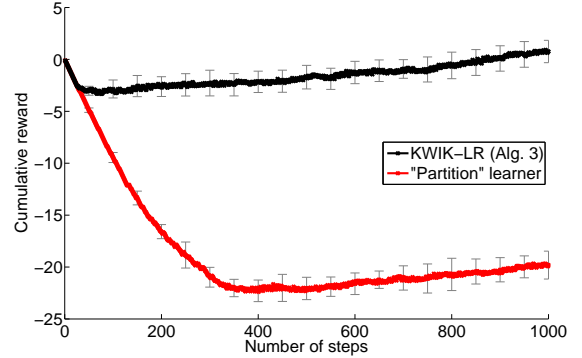


Figure 4: Stochastic OORL results for KWIK-LR and Partition in the Maze domain.

$c_i \in \mathcal{C}$ is satisfied (if any), and applies the corresponding effect $w_i \in \Omega^a$, which updates the attributes of the affected objects. Stochastic OOMDP effects generalize this representation so that a given condition induces not just a single effect (ω) but a distribution over possible effects (Ω^{ca}), only one of which (ω_i) will actually occur on a single timestep. In the parlance of action schemas as defined above, each $\{c, \Omega\}$ pair for action a can be thought of as its own action with preconditions specified by c . Again, this model can be viewed as a specific version of the action schemas presented above.

Previous work [5] presented an efficient algorithm for learning deterministic effects. Here, we demonstrate Algorithm 3 learning the probabilities associated with each effect in the stochastic setting when the possible effects given a condition and the conditions themselves are known in advance. Methods for relaxing these assumptions are discussed in later sections.

We demonstrate Algorithm 3 for stochastic OOMDPs on a simple 5×5 Maze domain, illustrated in Figure 3. The agent starts at location S and the goal is to arrive at G . Each step has a cost of -0.01 , and arriving at the goal results in a reward of $+1$. The agent’s actions are N, S, E and W. When executing an action, the agent will attempt to move in the desired direction with probability 0.8 and will slip to either side with probability 0.1. If it hits a wall, it stays put. This rule is what produces ambiguity in the effects. For example, imagine the agent has a wall to its North and East. If it attempts the N action, it could move to the West (with probability 0.1), or stay in place. If it stays in place, it might be because it attempted to move North (with probability 0.8) and hit the North wall, or it attempted to move East (with probability 0.1) and hit the East wall.

Figure 4 presents the results. As in the previous experiment, we can see that the KWIK-RL algorithm learns

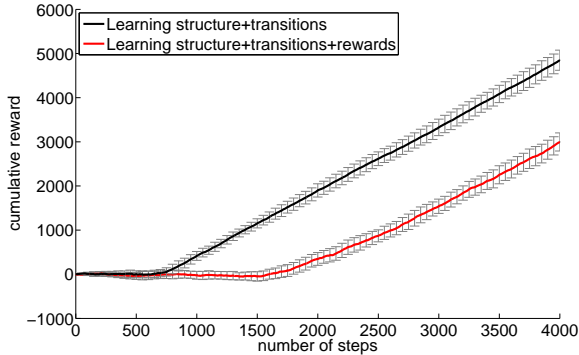


Figure 5: Average cumulative rewards over 10 runs on the Stocks domain.

much faster than the “Partition” learner by sharing information between equivalence partitions.

7 Extensions

We now discuss extensions for learning more of the compact representations discussed in this work. We will outline how to learn full fMDPs and larger parts of STRIPS and OOMDP models by combining KWIK-LR with other KWIK learners, and provide empirical support in the fMDP case. We also discuss a variation of KWIK-LR that can be used to learn stochastic action-schema outcomes (Ω^a).

7.1 Combining with Other KWIK Learners

Existing work [11] describes methods for combining simple KWIK agents to learn in increasingly complex domains. We follow this “building blocks” approach and show how to combine the fMDP reward-learning algorithm (Algorithm 2) with an implementation of the Noisy-Union [11; 10] algorithm (which is also KWIK) to learn the transition structure (Γ), transition probabilities (P), and reward function (R) of an fMDP all at once. The only knowledge given to the agent is the number of parents a factor may have ($|Z|$) and the reward-function structure, resulting in the only algorithm to date that learns all of these parameters efficiently. Experience is fed to both building block algorithms in parallel. The reward learner outputs an optimistic approximation of the reward function, which is given to Noisy-Union, which then learns the transition structure and probabilities. We conducted experiments to validate this algorithm in the Stocks domain. For comparison, we also ran the Noisy-Union algorithm with the rewards given *a priori*. Figure 5 displays the results, which show that all three quantities can be learned with small overhead.

We note that there are many other settings that could benefit from combining KWIK-LR with other KWIK learners for different parts of the model. For instance, in stochastic STRIPS and OOMDPs, the preconditions of actions (STRIPS) or the conditional effects (OOMDP) can be learned using an existing KWIK adaptation of Noisy Union [6] as long as their size is bounded by some known constant. Together, these learners could be used to learn all but the effects (Ω^a) of each action operator, a problem we now consider.

7.2 Future Work: Learning Effects

All of the effect distribution / condition learning variations require the possible effects (Ω^a) as input. Unfortunately, relaxing this assumption in the stochastic case is unlikely, since the effect learning problem is known to be NP-Hard [12]. When the number of possible effects is very small, or each ω_i is of constant size, enumeration techniques could be used. But, these assumptions are often violated, so researchers have concentrated on heuristic solutions [12]. Here, we suggest a novel heuristic that extends KWIK-LR for probability learning to the setting where the whole action schema (including Ω^a) needs to be learned.

We propose using sparsification; we consider all possible effects and use KWIK-LR to learn their probabilities, with an extra constraint that the number of “active” (non-zero) probabilities should be small. This minimization can be efficiently computed via linear programming, and techniques such as column generation may be used to keep the number of active constraints small in most cases. Together with the other learners discussed in this paper, the implementation of such an extension would form a complete solution to the action-schema learning problem: the probabilities and preconditions can be learned with KWIK-learners and the effects themselves can be learned heuristically using this sparsification extension. It remains a matter for future work to compare this system to other heuristic solutions [12] for such problems.

8 Related Work

Online linear regression has also been studied in the regret minimization framework (see e.g. [1]). Applications to restricted RL problems also exist [2], but with different types of bounds. Furthermore, regret analysis seems to lack the modularity (the ability to combine different learners) of the KWIK framework.

Previous work on linear regression for model-based RL has focussed on learning linear transition functions in continuous spaces. However, these approaches often lacked theoretical guarantees or placed restrictions on

the environment regarding noise and the reward structure [7]. In this paper, we have both improved on the current state-of-the-art algorithm for linear regression in RL [16], and used it in applications beyond the standard linear transition models. Our theoretical results and experiments illustrate the potential of KWIK-LR in model-based RL. In the future, we intend to identify other compact models where this technique can facilitate efficient learning and perform expanded empirical and theoretical studies of the extensions mentioned above.

Acknowledgements

We thank Lihong Li for spurring the discussion and that led to the revised analysis in this Technical Report. We also thank Csaba Szepesvári, Alexander L. Strehl, and the UAI reviewers for helpful comments. Support was provided by the Fulbright Foundation, DARPA IPTO FA8750-05-2-0249, FA8650-06-C-7606, and NSF IIS-0713435.

References

- [1] P. Auer. An improved on-line algorithm for learning linear evaluation functions. In *COLT*, 2000.
- [2] P. Auer, N. Cesa-Bianchi, and C. Gentile. Adaptive and self-confident on-line learning algorithms. *JCSS*, 64(1):48–75, 2002.
- [3] J. A. Boyan. Least-squares temporal difference learning. *Machine Learning*, 49(2):233–246, November 2002.
- [4] R. I. Brafman and M. Tennenholtz. R-MAX—a general polynomial time algorithm for near-optimal reinforcement learning. *JMLR*, 3:213–231, 2002.
- [5] C. Diuk, A. Cohen, and M. L. Littman. An object-oriented representation for efficient reinforcement learning. In *ICML*, 2008.
- [6] C. Diuk, L. Li, and B. Leffler. The adaptive k-meteorologists problem and its application to structure learning and feature selection in reinforcement learning. In *ICML*, 2009.
- [7] C.-N. Fiechter. PAC adaptive control of linear systems. In *COLT*, pages 72–80, New York, NY, USA, 1997. ACM.
- [8] R. Fikes and N. J. Nilsson. STRIPS: A new approach to the application of theorem proving to problem solving. *Artificial Intelligence*, 5:189–208, 1971.
- [9] M. J. Kearns and D. Koller. Efficient reinforcement learning in factored MDPs. In *IJCAI*, pages 740–747, 1999.
- [10] L. Li. *A Unifying Framework for Computational Reinforcement Learning Theory*. PhD thesis, Rutgers University, Piscataway, NJ, USA, 2009.
- [11] L. Li, M. L. Littman, and T. J. Walsh. Knows what it knows: A framework for self-aware learning. In *ICML*, 2008.
- [12] H. M. Pasula, L. S. Zettlemoyer, and L. P. Kaelbling. Learning symbolic models of stochastic domains. *JAIR*, 29:309–352, 2007.
- [13] K. B. Petersen and M. S. Pedersen. The matrix cookbook, 2005. Version 20051003.
- [14] A. A. Sherstov and P. Stone. Improving action selection in MDP’s via knowledge transfer. In *AAAI*, 2005.
- [15] A. L. Strehl, C. Diuk, and M. L. Littman. Efficient structure learning in factored-state MDPs. In *AAAI*, 2007.
- [16] A. L. Strehl and M. L. Littman. Online linear regression and its application to model-based reinforcement learning. In *NIPS*, 2008.
- [17] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, March 1998.
- [18] I. Szita and A. Lőrincz. The many faces of optimism: a unifying approach. In *ICML*, pages 1048–1055, 2008.
- [19] S. A. van de Geer. *Applications of Empirical Process Theory*. Cambridge University Press, 2000.
- [20] H. L. S. Younes, M. L. Littman, D. Weissman, and J. Asmuth. The First Probabilistic Track of the International Planning Competition. *JAIR*, 24:851–887, 2005.