

Reinforcement Planning: RL for Optimal Planners

Matt Zucker and J. Andrew Bagnell

CMU-RI-TR-10-14

April 2010

Robotics Institute
Carnegie Mellon University
Pittsburgh, Pennsylvania 15213

© Carnegie Mellon University

Abstract

Search based planners such as A* and Dijkstra’s algorithm are proven methods for guiding today’s robotic systems. Although such planners are typically based upon a coarse approximation of reality, they are nonetheless valuable due to their ability to reason about the future, and to generalize to previously unseen scenarios. However, encoding the desired behavior of a system into the underlying cost function used by the planner can be a tedious and error-prone task. We introduce *Reinforcement Planning*, which extends gradient based reinforcement learning algorithms to automatically learn useful cost functions for optimal planners. Reinforcement Planning presents several advantages over other learning applications involving planners in that it is not limited by the expertise of a human demonstrator, and that it also recognizes that the domain of the planner is a simplified model of the world. We demonstrate the effectiveness of our method in learning to solve a noisy physical simulation of the well-known “marble maze” toy.

Contents

1	Introduction	1
1.1	Background and Related Work	1
1.2	Outline	3
2	Value Function Subgradient	3
3	Reinforcement Planning	4
3.1	Value function approximation	5
3.2	Direct policy learning	6
3.3	Discussion of RP implementations	7
4	Example Domain: Marble Maze	7
4.1	Marble maze dynamics and reward	8
4.2	Marble maze planner and policy	9
5	Results	10
6	Conclusion	12

1 Introduction

State-of-the-art robotic systems [1, 2, 3] increasingly rely on search-based planning or optimal control methods to guide decision making. Similar observations can be made about computer game engines. Such methods are nearly always extremely crude approximations to the reality encountered by the robot: they consider a simplified model of the robot (as a point, or a “flying brick”), they often model the world deterministically, and they nearly always optimize a *surrogate* cost function chosen to induce the correct behavior rather than the “true” reward function corresponding to a declarative task description. Such approximations are made as they enable efficient, real-time solutions.

Despite this crudeness, optimal control methods have proven quite valuable because of the ability to transfer knowledge to new domains; given a way to map features of the world to a cost function, we can compute a plan that navigates a robot in a never-before-visited part of the world. While value-function methods and reactive policies are popular in the reinforcement learning (RL) community, it often proves remarkably difficult to transfer the ability to solve a particular problem to related ones using such methods [4]. Planning methods, by contrast, consider a sequence of decisions in the future, and rely on the principle underlying optimal control **that cost functions are more parsimonious and generalizable than plans or values.**

However, planners are only successful to the extent that they can transfer domain knowledge to novel situations. Most of the human effort involved in getting systems to work with planners stems from the tedious and error-prone task of adjusting surrogate cost functions, which has until recently been a black art. Imitation learning by Inverse Optimal Control, using, *e.g.* the Learning to Search approach [3], has proven to be an effective method for automating this adjustment; however, it is limited by human expertise.

Our approach, Reinforcement Planning (RP), is straightforward: we demonstrate that crude planning algorithms can be learned as part of a direct policy search or value function approximator, thereby allowing a system to experiment on its own and outperform human expert demonstration.

1.1 Background and Related Work

Central to this work is the distinction between the “true” reward function for a robotic system, and the cost function used by an optimal planner. Consider an unmanned ground vehicle (UGV) system similar to that discussed in [5]. Modern UGV software architectures are typically hierarchical planning and control schemes, as depicted in Figure 1: a coarse global planner computes a path from the robot’s current position to a distant goal, and a “local planner” or control policy runs at a much higher frequency in order to guide the vehicle along the planned path. The coarse planner for our hypothetical vehicle may neglect to model the kinematics and dynamics of the underlying robotic platform, reducing the planning problem to motion on an 8-way-connected, regularly sampled grid.

The true reward function for the UGV control system is simple to write down: it encompasses some assessment of the current robot state x , evaluated at every control cycle until termination:

$$r(x) = \begin{cases} \$100, & \text{if goal is reached} \\ -\$100,000, & \text{if vehicle capsizes} \\ -\$10,000, & \text{if hit an immovable obstacle} \\ \vdots & \\ -\$0.01, & \text{otherwise} \end{cases} \quad (1)$$

However, because the planner reasons over a crude approximation of reality, it is unable to optimize the true reward function directly. Instead, the planner reasons about a surrogate cost function, which might consist

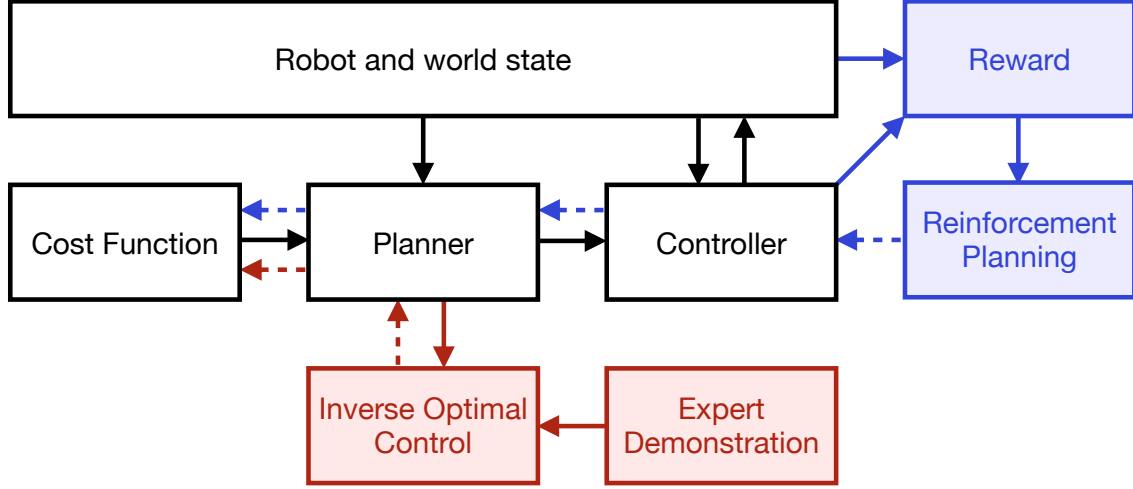


Figure 1: A schematic comparison of Reinforcement Planning and Inverse Optimal Control. Solid arrows denote information flow and causal relationships; dotted arrows indicate feedback from learning algorithms. The black outlined boxes in the upper left illustrate typical components of a hierarchical planning system. The red components below constitute the Inverse Optimal Control approach to learning cost functions, and the blue components to the right correspond to our novel Reinforcement Planning method.

of a linear combination of scalar features associated with a discrete grid state s :

$$c(s) = \theta^T f(s), \quad f(s) = \begin{pmatrix} \text{slope} \\ \text{roughness} \\ \text{visibility} \\ \text{traversability} \\ \text{clutter} \\ \vdots \end{pmatrix} \quad (2)$$

This disconnect between true reward and surrogate cost is precisely what makes tuning cost functions difficult for system designers; their task is to identify a cost function which implicitly causes the system to behave favorably in terms of the actual reward.

The Learning to Search (LEARCH) approach [3] removes the burden of manually parameterizing cost functions by allowing a system designer to specify the desired output of the coarse planner. Whereas the goal of optimal control is to identify a path or trajectory that minimizes a cost function, *Inverse Optimal Control* (IOC) schemes such as LEARCH attempt to find a cost function under which demonstrated paths or trajectories appear optimal.

Although LEARCH is powerful due to its convex formulation and its ability to reduce development time for planning systems, the approach has limitations. First and foremost, it is limited by the performance of the expert demonstration. A human may be unskilled at the target task, or worse, unable to demonstrate it at all. Furthermore, LEARCH is also fundamentally limited by the fact that its learning feedback does not generally consider the local planner or control policy: in some sense, it is learning at the wrong level of abstraction. Our Reinforcement Planning method, on the other hand, propagates the reward signal back

through the controller and planner to the parameters of the underlying cost function, and is not limited by any expert demonstration.

We also note that planning algorithms have been used extensively in previous RL work (*e.g.* Dyna [6]); our work contrasts with these in embracing the reality that real-world planners operate on a coarse approximation and must be trained to have a cost-function that induces the correct behavior, not the “true” cost function, and so additional indirection is required to solve the RL problem. Moreover, approaches such as Dyna attempt to (approximately) solve the full Markov Decision Problem (MDP) associated with the system. For the types of complex high-dimensional systems that are considered by hierarchical planning approaches, the curse of dimensionality makes an explicit representation of a full value function or policy corresponding to the solution of such an MDP computationally intractable.

1.2 Outline

The remainder of this document is organized as follows: in [Section 2](#), we outline the key observation that underpins RP, the subgradient of a planner’s value function. In [Section 3](#), we detail how RP can be implemented in terms of the value function subgradient, both for direct policy gradient algorithms, and for value function approximation schemes. In [Section 4](#), we illustrate RP on a practical example: solving a noisy physical simulation of the marble maze. We review the results of our simulation experiments in [Section 5](#), and finally, we conclude in [Section 6](#).

2 Value Function Subgradient

The central observation underlying RP is that the value function computed by an optimal planner has a subgradient. If a policy for the system is computed in terms of the planner’s value function, the subgradient of the value function will appear in the policy gradient wherever the value function appears in the policy. Subsequently, any gradient based reinforcement learning algorithm can be straightforwardly extended to consider optimal planners.

The task for a discrete optimal planner is to find the minimum-cost sequence of actions to transition from a starting state to a goal state. Let $s \in \mathcal{S}$ denote a discrete state, and let $a \in \mathcal{A}$ denote an action. The successor state given an action is specified by a state transition model $s' = \text{succ}(s, a) : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$. Denote by $\xi(s_0)$ a sequence of state-action pairs (s_k, a_k) which starts at s_0 and reaches a goal state in $\mathcal{S}_{\text{goal}} \subset \mathcal{S}$. Given a one-step cost function $c(s, a) : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$, the formal task for the planner is to compute $\xi^*(s_0)$, a minimum-cost path from s_0 to the goal. The *value* $V(s_0)$ of the state s_0 is defined as the sum of costs along that path. The *value function* $V(s)$ is then the function which maps each state s to the corresponding minimum cost:

$$V(s) = \min_{\xi(s)} \sum_{(s_k, a_k) \in \xi(s)} c(s_k, a_k) \quad (3)$$

Although we traditionally consider algorithms such as A*, and D* [7, 8] to be chiefly concerned with identifying $\xi^*(s)$, a path which minimizes $V(s)$ above, they can equivalently be considered to compute $V(s)$ itself.¹ There is a computational tradeoff between using “path-oriented” algorithms such as A* and D* and using “value-oriented” algorithms such as Dijkstra’s [9] which compute the value of all states. For

¹Either problem may be trivially transformed into the other: given an optimal trajectory, the corresponding value is simply the sum of costs along that trajectory; given a value function, an optimal trajectory can be computed via steepest descent.

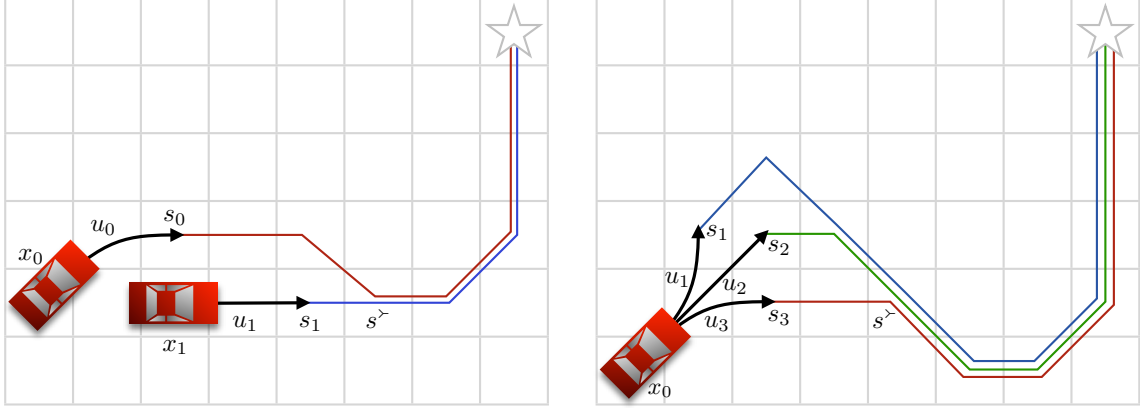


Figure 2: Illustration of SARSA(λ)-style (*left*) and REINFORCE-style (*right*) gradients under RP. In both cases, gradient terms involve a sum of features along prefixes of optimal paths up to s^λ . See discussion in [Section 3.1](#) and [Section 3.2](#) for details.

many scenarios, precomputing a full value function over the state space may waste computation if the set of desired values is far smaller than the full state space.

Suppose that the cost function is defined as a linear combination of features defined over state-action pairs, based on a weighting θ :

$$c(s, a, \theta) = \theta^T f(s, a) \quad (4)$$

Then the subgradient of $V(s, \theta)$, as written in [Equation 3](#), with respect to θ is simply the sum of the one-step-cost gradients along the optimal path $\xi^*(s)$, namely the features themselves:

$$\nabla V(s, \theta) = \sum_{(s_k, a_k) \in \xi^*(s)} \nabla c(s_k, a_k, \theta) \quad (5)$$

$$= \sum_{(s_k, a_k) \in \xi^*(s)} f(s_k, a_k) \quad (6)$$

This “feature counting” interpretation of differential value has been described before in the derivation of Maximum Margin Planning and other LEARCH algorithms [\[3\]](#); however, instead of directly matching feature counts with those of an expert as in LEARCH, our goal in Reinforcement Planning is to provide a way for a learning system to modify planner values in order to maximize expected reward.

3 Reinforcement Planning

Let $x \in \mathcal{X}$ be the full state (or observation) space for a robotic system, and let $u \in \mathcal{U}$ be the space of controls. We assume the existence of a function $\text{proj}(x) : \mathcal{X} \mapsto \mathcal{S}$ which maps a full state x to the corresponding “nearby” state s in the domain of the coarse planner. An admissible *Q-function* for Reinforcement Planning is a function $Q(x, u, \theta) \mapsto \mathbb{R}$ which computes the relative quality of applying action u at state x in terms of the optimal cost-to-go $V(s', \theta)$ for some coarse planner state(s) s' that depends on x and u .

For example, assume there is a deterministic state transition function $x' = \text{succ}(x, u)$ for the system. Then one admissible Q-function for RP would be to compute the optimal cost-to-go of the projected successor state:

$$Q(x, u, \theta) = V(\text{proj}(\text{succ}(x, u)), \theta) \quad (7)$$

If, as in a traditional MDP, there is no deterministic state transition function, we would instead consider the expected cost-to-go under the distribution $p(x'|x, u)$:

$$Q(x, u, \theta) = \sum_{x'} V(\text{proj}(x'), \theta) p(x'|x, u) \quad (8)$$

We now proceed by constructing a policy (in the RL sense) around Q , and adapting existing gradient based RL algorithms to optimize θ based on measured rewards. Any such gradient based RL algorithm will examine ∇Q , the gradient of Q with respect to the cost function parameters. We note that anywhere V appears in the definition of Q , so too will ∇V appear in the expression of ∇Q .²

Broadly speaking, the taxonomy of gradient based RL algorithms can be broken down into those which use value function approximators, and those which learn policies directly. The following subsections illustrate how to apply RP in both cases.

3.1 Value function approximation

The Q-function defined above can be considered as a parametric value function approximator. By constructing an ε -greedy policy around Q , we can learn the cost function parameters θ via steepest descent on Bellman Error with an eligibility trace, very similar to the derivation of SARSA(λ) [10].

To begin with, we construct an ε -greedy policy on Q . Let $U(x) \subseteq \mathcal{U}$ be the set of actions available at state x . Then define the policy $\pi(x, \theta)$ as

$$\pi(x, \theta) = \begin{cases} \arg \min_{u \in U(x)} Q(x, u, \theta), & \text{with probability } (1 - \varepsilon) \\ \text{random } u \in U(x), & \text{with probability } \varepsilon \end{cases} \quad (9)$$

Subsequently, we follow the RP-VFA algorithm, defined below.

²Practically, other terms may appear in the right hand side of Equation 7 and Equation 8. For instance, if u is a vector of actuator torques, adding $\kappa \|u\|^2$ (with κ a scalar weight) is a typical way to minimize energy while pursuing a goal. In general, these terms will disappear when computing the gradient of Q with respect to θ ; however, it may be desirable to optimize parameters such as κ in tandem.

Algorithm 1: RP-VFA algorithm for steepest descent on Bellman Error with eligibility trace.

```

1 initialize  $e(x, u) \leftarrow 0$  for all  $x, u$ ;
2 initialize  $\Delta \leftarrow 0$ ;
3 pick initial  $x, u$ ;
4 while  $x$  is not terminal do
5   take action  $u$ , observe reward  $r(x, u)$  and successor state  $x'$ ;
6   choose  $u'$  from  $\pi(x', \theta)$ ;
7    $\delta \leftarrow Q(x, u, \theta) - Q(x', u', \theta) - r(x, u)$ ;
8    $e(x, u) \leftarrow e(x, u) + \nabla Q(x, u, \theta) - \nabla Q(x', u', \theta)$ ;
9   foreach previous  $x, u$  do
10     $\Delta \leftarrow \Delta + \delta e(x, u)$ ;
11     $e(x, u) \leftarrow \lambda e(x, u)$ ;
12   end
13    $x \leftarrow x'; u \leftarrow u'$ ;
14 end
15  $\theta \leftarrow \theta + \alpha \Delta$ ;

```

We note that the definition of $e(x, u)$ on [line 8](#) is the gradient of the Bellman Error δ with respect to the parameters θ , accumulated through the eligibility trace. The left hand side of [Figure 2](#) illustrates the effect of the gradient term: initially, the system is at state x_0 and chooses action u_0 based on the value of the underlying coarse planner state s_0 . After executing action u_0 , the system lands in state x_1 , incurring reward $r_0 = r(x_0, u_0)$. The next action chosen is u_1 , based on the value of the projected successor state s_1 . The gradient terms computed in [line 8](#) above are equal to the difference of the features encountered along the two optimal paths in the figure. Note that both optimal paths to the goal merge at the state s^* ; hence any features along the shared subpath from s^* to the goal cancel out and do not appear in the resulting value of $e(x, u)$.

The practical effect of λ is to smear the impact of rewards backwards in time, which is vital for learning under the presence of uncertainty or modeling artifacts due to the planner’s coarse approximation of reality. After some number of training episodes, Δ holds the overall gradient direction, along which we move θ by some step size α .

3.2 Direct policy learning

Instead of considering Q as a value function approximator, we can use it as the basis for a stochastic policy, and subsequently draw upon direct policy learning algorithms such as REINFORCE [[11](#)] to learn θ . We begin by constructing a stochastic policy $p(u|x, \theta)$ around Q :

$$p(u_i|x, \theta) = \frac{1}{Z} \exp(-\nu Q(x, u_i, \theta)) \quad (10)$$

where Z is the normalizer that causes $p(u|x, \theta)$ to sum to one over all $u \in U(x)$, and ν is positive scalar weight. REINFORCE requires us to compute the *score ratio* $\nabla p_i/p_i$. For the Boltzmann-style distribution defined in [Equation 10](#), the score ratio can be expressed as

$$\frac{\nabla p(u_i|x, \theta)}{p(u_i|x, \theta)} = -\nu (\nabla Q(x, u_i, \theta) - E[\nabla Q(x, u, \theta)]) \quad (11)$$

$$= -\nu \left(\nabla Q(x, u_i, \theta) - \sum_{u \in U(x)} p(u|x, \theta) \nabla Q(x, u, \theta) \right) \quad (12)$$

At this point, the REINFORCE algorithm can be applied straightforwardly. We reproduce it here, following the derivation in [12]:

Algorithm 2: REINFORCE for direct policy learning with RP

```

1 initialize  $z_0 \leftarrow 0, \Delta_0 \leftarrow 0$ ;
2 initialize  $t \leftarrow 0$ ;
3 initialize  $x_0$ ;
4 while  $x_t$  is not terminal do
5   sample  $u_t$  from  $p(u_t|x_t, \theta)$ ;
6   execute  $u_t$ , observe reward  $r_t \leftarrow r(x_t, u_t)$  and successor state  $x_{t+1}$ ;
7    $z_{t+1} \leftarrow \beta z_t + \frac{\nabla p(u_t|x_t, \theta)}{p(u_t|x_t, \theta)}$ ;
8    $\Delta_{t+1} \leftarrow \Delta_t + \frac{1}{t+1} (r_t z_t - \Delta_t)$ ;
9    $t \leftarrow t + 1$ ;
10 end
11  $\theta \leftarrow \theta + \alpha \Delta_t$ ;

```

As with the value function approximation algorithm, α is a step size parameter. Here, $\beta \in [0, 1]$ is a discount factor. The score ratio for RP is illustrated on the right hand side of Figure 2: consider the robot at state x_0 . The respective probabilities p_i of the actions u_i are based on the optimal cost-to-go of the projected successor states s_i , for $i \in 1, 2, 3$. Say that we sample action u_1 . Then the score ratio $\nabla p_1/p_1$ is defined to be the difference between the sum of features along the blue path from s_1 , and the expected features over all three paths. As with the value function approximation example, all optimal paths merge at a future planner state s^\star , and therefore the features along the common subpaths are eliminated from the computation of the score ratio in Equation 12.

3.3 Discussion of RP implementations

As illustrated in Figure 2, the gradient update rules for both approaches modify the parameters θ based upon observed rewards, and a set of features along optimal paths to the goal. However, in general, neither algorithm counts features all the way to the goal state: in both cases, there is typically a state s^\star at which the sampled optimal paths combine, causing features along the common subpaths to cancel out. One interpretation of the RP gradient terms is to consider them as a comparison of features which differ among various optimal paths found by the low-level planner.

As a gradient based method, Reinforcement Planning is by no means guaranteed to converge to a global optimum of θ for either type of approach. Therefore, initializing θ to a principled estimate is necessary in practice. If *a priori* estimates of θ are not available, using LEARCH or a similar IOC approach is a suitable initialization method.

4 Example Domain: Marble Maze

The familiar “marble maze” or labyrinth toy (shown in Figure 3) has been used in the past as a benchmark application for learning algorithms, particularly imitation learning [4, 13]. However, previous approaches learned policies specifically targeted at a particular board geometry, and failed to generalize to novel unseen boards without significant additional learning [14]. Just as more complex dynamical systems such as legged robots, the marble maze task exhibits significant momentum effects, and requires reasoning about the

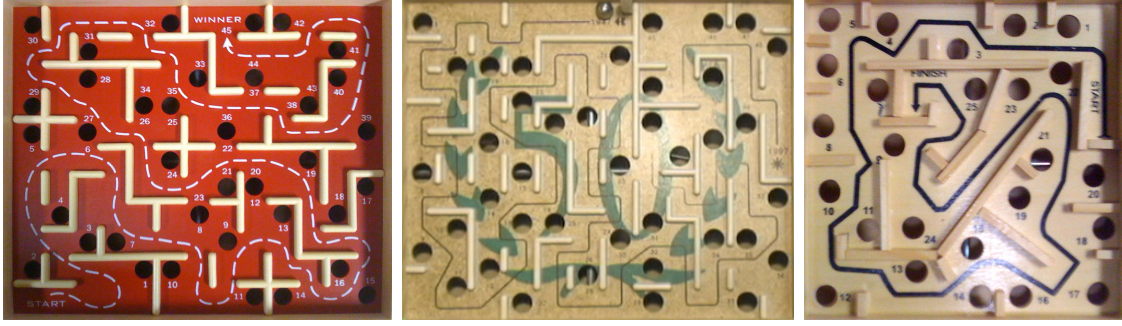


Figure 3: A representative sampling of marble mazes. The operator must tilt the board using two knobs (not shown) to guide a marble from the starting position to the goal position.

consequences of future actions. We chose to demonstrate Reinforcement Planning on the marble maze task in order to show that a very simple planner can produce policies that not only perform well in the face of dynamics and uncertainty, but also generalize to novel mazes.

The marble maze consists of a flat board with raised walls. An operator may apply forces to the marble by turning knobs which rotate the board along its x and y axes. The goal of the game is to successfully navigate the marble from a starting position to a goal position, avoiding any number of holes drilled into the board. Falling into a hole ends the game.

4.1 Marble maze dynamics and reward

We define the state of the marble maze as

$$x = (p_x, p_y, \dot{p}_x, \dot{p}_y, r_x, r_y)^T$$

which combines the 2D position p_x, p_y and velocity \dot{p}_x, \dot{p}_y of the ball with the current tilt angles r_x, r_y of the board. The command for the system $u = (u_x, u_y)^T$ corresponds to desired tilt angles for the board. The reward for the system depends only upon the state. Evaluated once per 10 Hz control update until the game has ended, it is defined as:

$$r(x) = \begin{cases} 100, & \text{if goal reached} \\ -100, & \text{if fell in hole} \\ -0.01 & \text{otherwise} \end{cases} \quad (13)$$

We model the marble maze in a physical simulation. The acceleration imparted on the ball by board tilt is given by:

$$\ddot{p}_x = g \sin(r_y) \quad (14)$$

$$\ddot{p}_y = -g \cos(r_y) \sin(r_x) \quad (15)$$

where $g = 9.8m/s^2$. There is additional acceleration due to rolling friction, which is computed to be proportional to the velocity of the ball. The coefficient of rolling friction is $\mu_{roll} = 0.05$. Collisions between the ball and a wall result in an instantaneous reflection of the velocity vector about the normal of the wall,

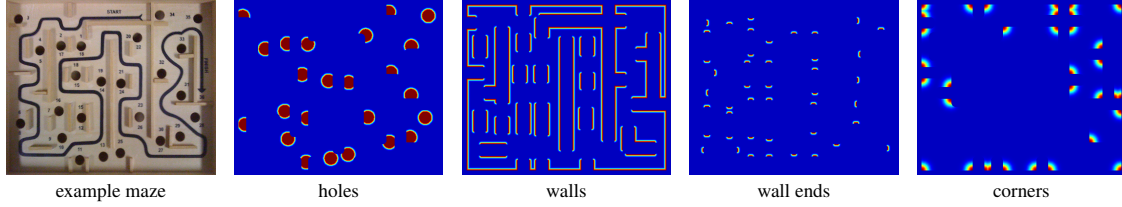


Figure 4: An example marble maze and features used for planning. All features are computed by “blurring” board geometry with two different kernel sizes (only one is shown here for each feature).

and damping by a coefficient of restitution $c_r = 0.85$. After the acceleration and impacts are computed, the position and velocity of the ball is updated via the midpoint method.

Second-order dynamics of the rotation of the board are ignored. We assume that the board tilt is highly position controlled, with

$$\dot{r}_x = \min \left(\frac{u_x - r_x}{\Delta t}, \dot{r}_{max} \right) \quad (16)$$

$$\dot{r}_y = \min \left(\frac{u_y - r_y}{\Delta t}, \dot{r}_{max} \right) \quad (17)$$

where Δt is the integration time step.

We also model a variety of adverse effects expected to be encountered in a true physical system such as control noise, local variation in μ_{roll} and c_r , and board warp (local variation in board tilt). Control noise is drawn from a normal distribution with mean zero and $\sigma_u = 0.02\text{rad}$. Local variations in the physical parameters are drawn from a spatially coherent noise distribution [15]. Finally, we also model differences between the idealized board considered by the planner and the underlying board used for simulation by slightly displacing walls and holes in the two models.

Our goal is to learn a policy mapping states to actions using RP. Again, as expected in a physical system, we present our policy with a noisy observation of state instead of the true simulated state of the system.

4.2 Marble maze planner and policy

As in Figure 1, we create a hierarchical planning architecture for the noisy marble maze simulation. We construct a coarse planner that operates over a 2D grid corresponding to the 2D workspace \mathcal{S} of the maze board. Each grid cell $s \in \mathcal{S}$ is assigned a set of features $f(s)$ representing aspects of the board geometry. Features, shown in Figure 4, are computed as blurred versions of holes, walls, wall ends, and corners. There are two different kernel radii used for the blurring. All feature values lie in the interval $[0, 1]$. We augment the feature vector with the logical complement of each feature $1 - f$, and add a constant bias feature of 1.0, for a total of 17 features in all.

We use Dijkstra’s algorithm [9] to compose a value function representing the optimal cost-to-go from each grid cell s to the goal. The one-step-cost is computed as $c(s) = \theta^T f(s)$ with θ a positive vector of weights. Grid cells lying inside holes and walls are treated specially by our value function algorithm: the value function in such location always points back to the nearest grid cell accessible to the ball.

The policy defined for the marble maze considers a discrete set of candidate actions. For each action u

at state x , we define the function

$$Q(x, u, \theta) = \frac{k_u}{2} \|u - u^*(x, \theta)\|^2 + k_v V(s', \theta) \quad (18)$$

Where k_u and k_v are scalar weights, $u^*(x, \theta)$ (defined below) is a commanded tilt that attempts to align the velocity of the ball with the local direction of the planner value function, and s' is the projected successor state $\text{proj}(\text{succ}(x, u))$. Projection is defined as finding the grid cell nearest to p_x, p_y .

We define $u^*(x, \theta)$ as a proportional controller on ball velocity with gain k_p . The desired velocity v_d comes from the normalized directional gradient of the planner value function $n(x, \theta)$ with respect to the board x - and y -axes, scaled by a gain ℓ . The error $\epsilon(x, \theta)$ between desired and current velocity is converted to a commanded tilt by a skew-symmetric matrix M which accounts for the accelerations in Equation 14 and Equation 15.

$$n(x, \theta) = \begin{bmatrix} \frac{\partial}{\partial p_x} \\ \frac{\partial}{\partial p_y} \end{bmatrix} V(\text{proj}(x), \theta) \quad (19)$$

$$v_d(x, \theta) = \ell \frac{n(x, \theta)}{\|n(x, \theta)\|} \quad (20)$$

$$\epsilon(x, \theta) = v_d(x, \theta) - (\dot{p}_x, \dot{p}_y)^T \quad (21)$$

$$M = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \quad (22)$$

$$u^*(x, \theta) = k_p M \epsilon(x, \theta) \quad (23)$$

To apply RP, we need to take the gradient of Equation 18 with respect to θ :

$$\nabla Q(x, u, \theta) = -k_u \nabla u^*(x, \theta) (u - u^*(x, \theta)) + k_v \nabla V(s', \theta) \quad (24)$$

where $\nabla u^*(x, \theta)$ is the 17-by-2 matrix

$$\nabla u^*(x, u, \theta) = k_p \ell \left(M \left(I - \frac{n n^T}{\|n\|^2} \right) \begin{bmatrix} \frac{\partial}{\partial p_x} \\ \frac{\partial}{\partial p_y} \end{bmatrix} \nabla V(\text{proj}(x), \theta)^T \right)^T \quad (25)$$

Above, I is the 2-by-2 identity matrix, and $n = n(x, \theta)$.

5 Results

We used RP with both the RP-VFA algorithm of Section 3.1 as well as the REINFORCE algorithm derived in Section 3.2 to learn a cost function for the coarse planner to solve the noisy marble maze simulation. The initial value for the weight vector θ was zero for all features, except for a small positive bias term which encouraged goal-seeking behavior. Reinforcement Planning was run on a very small training board for 100 trial episodes with both learning methods, and the performance was evaluated on the training board as well as two significantly larger test boards (shown in Figure 5).

We compared RP to a simple regression algorithm that attempts to match the one-step-costs of the planner to the measured rewards in the simulator. The regression algorithm failed to learn any significant structure of the problem, presumably due to the fact that it was unable to ascribe any of the current reward to features

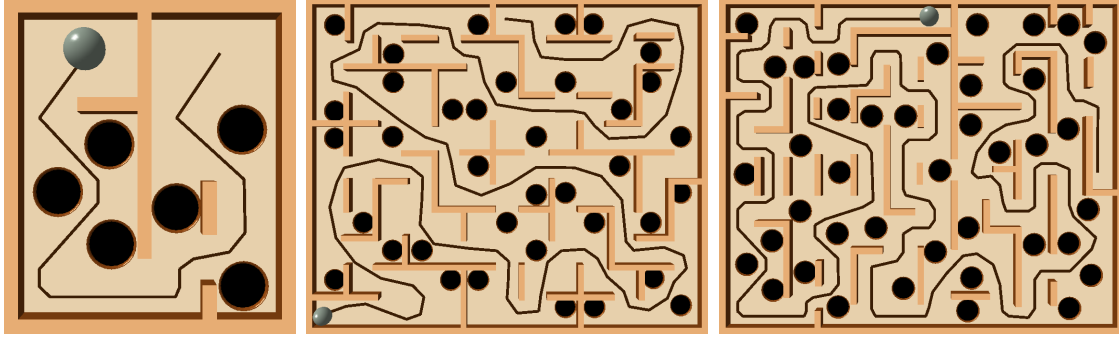


Figure 5: *Left to right:* Training board, test board 1, and test board 2. The test boards are successfully navigated by a policy learned with RP despite being trained only upon the small test board.

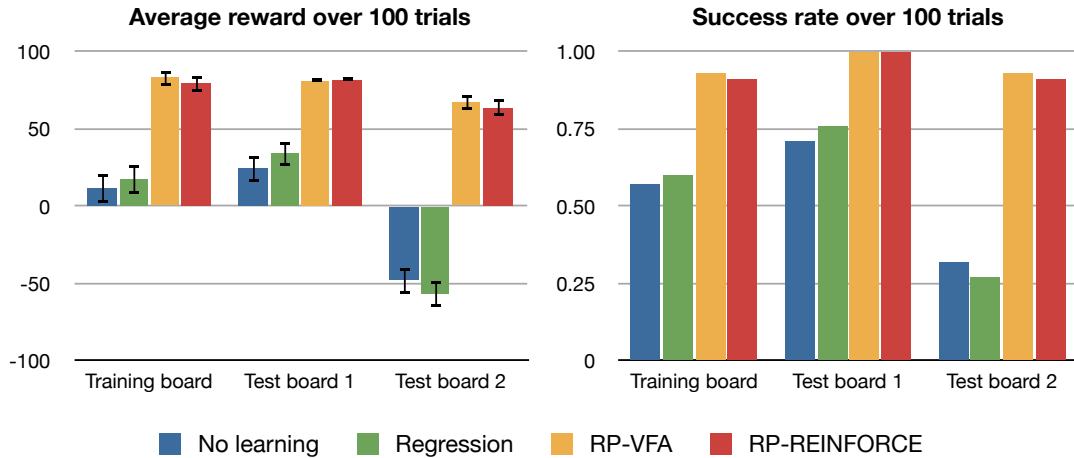


Figure 6: Results of experiments on noisy marble maze simulation after training for 100 episodes.

previously encountered by the system. This “smearing” back in time of the correspondence between rewards and features is handled well by both the REINFORCE and the RP-VFA algorithms.

Reinforcement Planning, however, was able to successfully learn how to solve not only the training board far better than the initial weights, but it also vastly outperformed the initial weights and the regression algorithm on the much larger, and previously unseen, test boards. The results are summarized in Figure 6. Learning converges after approximately 100 trials; training for many more trials did not significantly alter performance.

The cost function learned by RP is shown in Figure 7. Reinforcement Planning confirms human experts’ own intuition about solving the marble maze problem: not only is cost high near holes, but it is lower next to walls and corners, presumably because their effect of reducing uncertainty of the rolling marble dynamics. Furthermore, the system has discovered the correct weighting of features on its own, obviating the need for hand-tuning cost function parameters. After normalizing the weight vectors to sum to one, it is clear the weights learned by the two algorithms are quite similar.

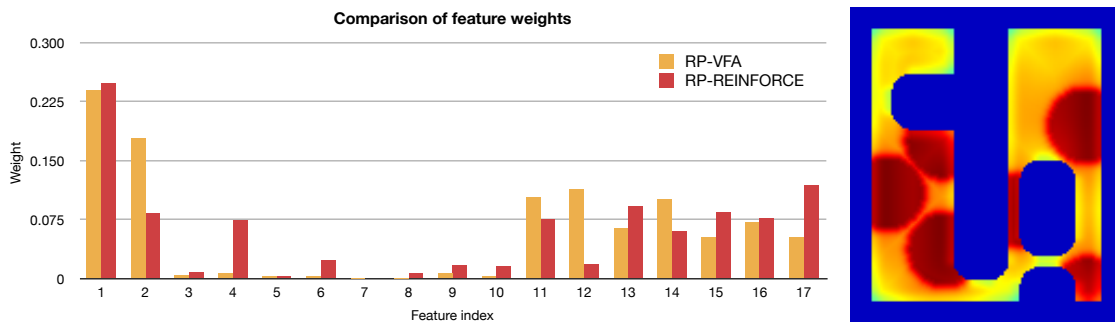


Figure 7: *Left*: comparison of feature weights learned by RP-VFA and RP with REINFORCE. Weight vectors have been normalized to sum to one. *Right*: Visualization of cost function learned by REINFORCE. RP not only learns that holes are associated with high cost, but that areas near walls and corners have lower cost than free space.

6 Conclusion

We have introduced Reinforcement Planning, a novel extension to existing reinforcement learning algorithms which allows them to reason about optimal planners. The key of RP is to define a policy for a learning system in terms of the value computed by a coarse planner. Computing the gradient of the policy evaluates the subgradient of the planner’s value function. Reinforcement Planning can be used in both value function approximation, and direct policy gradient settings. We experimentally verified the effectiveness of RP in learning to control a noisy physical simulation of the marble maze game.

In future work, we will investigate a number of extensions to the methods proposed here, and begin to implement them on physical robotic systems as opposed to computer simulations. We believe it will be straightforward to derive the functional analogues of the algorithms put forward in [Section 3.1](#) and [Section 3.2](#): just as MMPBoost [16] is a boosted version of Maximum Margin Planning [3], so too can we derive boosted versions of RP for automatic feature discovery or non-linear cost functions.

Our derivations of Reinforcement Planning in this paper have focused on discrete search as the underlying implementation of optimal planners. However, other optimal planning and control schemes exist, such as variational methods which use Pontryagin’s minimum principle or Differential Dynamic Programming [17], which operates on continuous state and action domains.

We have begun work on a physical implementation of the marble maze platform (shown in [Figure 8](#)), using hardware courtesy of Martin Stolle [4]. In future work, we hope to demonstrate a solution of the marble maze on the physical robot platform. We would also like to investigate learning cost functions for the LittleDog quadruped footstep planner [18].

References

- [1] Joel Chestnutt. *Navigation Planning for Legged Robots*. PhD thesis, Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, December 2007.
- [2] C. Urmson et al. Autonomous driving in urban environments: Boss and the urban challenge. *Journal of Field Robotics*, 2008.

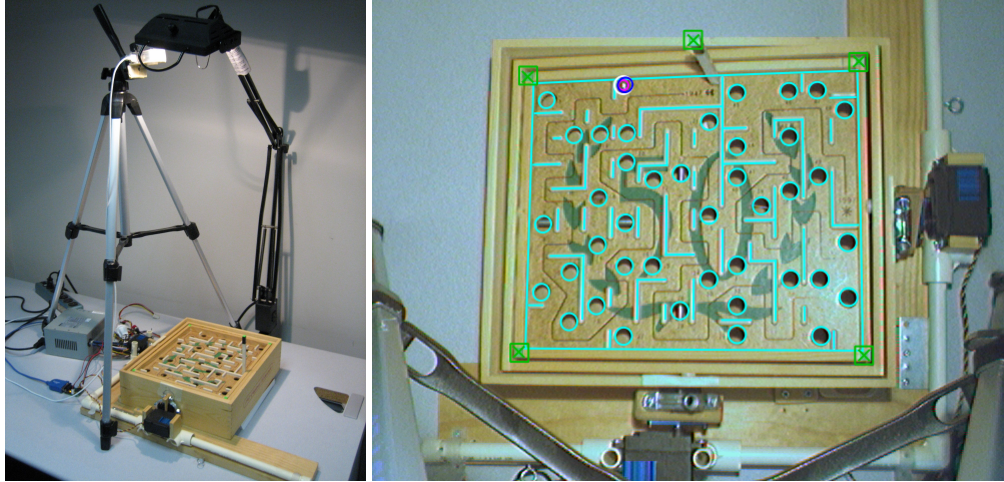


Figure 8: Physical robotic platform for playing the marble maze. *Left*: the robot consists of two hobby servos, and an overhead camera. *Right*: a computer vision system recognizes the position of the ball and markers on the board to reconstruct the 3D pose of the system.

- [3] N. Ratliff, D. Silver, and J. A. Bagnell. Learning to search: Functional gradient techniques for imitation learning. *Autonomous Robots*, 2009.
- [4] Martin Stolle. *Finding and Transferring Policies Using Stored Behaviors*. PhD thesis, Robotics Institute, Carnegie Mellon University, Pittsburgh, PA, 2008.
- [5] D. Silver, J. A. Bagnell, and A. Stentz. High performance outdoor navigation from overhead data using imitation learning. In *Robotics Science and Systems*, June 2008.
- [6] R. S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *ICML*, 1990.
- [7] P.E. Hart, N.J. Nilsson, and B. Raphael. A Formal Basis for the Heuristic Determination of Minimum Cost Paths. *Systems Science and Cybernetics, IEEE Transactions on*, 4(2):100–107, 1968.
- [8] A. Stentz. The focussed D* algorithm for real-time replanning. *Proceedings of the International Joint Conference on Artificial Intelligence*, pages 1652–1659, 1995.
- [9] E.W. Dijkstra. A note on two problems in connexion with graphs. *Numerische mathematik*, 1(1):269–271, 1959.
- [10] R.S. Sutton and A.G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [11] Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 1992.
- [12] J. Baxter and P.L. Bartlett. Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 2001.

- [13] Darrin Bentivegna. *Learning from Observation Using Primitives*. PhD thesis, Georgia Institute of Technology, 2004.
- [14] Martin Stolle and Chris Atkeson. Knowledge transfer using local features. In *Proceedings of the IEEE Symposium on Approximate Dynamic Programming and Reinforcement Learning*, 2007.
- [15] K. Perlin. Improving noise. *ACM Transactions on Graphics (TOG)*, 2002.
- [16] Nathan Ratliff, David Bradley, J. Andrew Bagnell, and Joel Chestnutt. Boosting structured prediction for imitation learning. In *NIPS*, Vancouver, B.C., December 2006.
- [17] D.H. Jacobson and D.Q. Mayne. *Differential dynamic programming*. Elsevier Publishing Company, 1970.
- [18] M. Zucker, J.A. Bagnell C.G. Atkeson, and J. Kuffner. An optimization approach to rough terrain locomotion. In *Proc. IEEE Int'l Conf. on Robotics and Automation*, 2010.