

Gradient-based Reinforcement Planning in Policy-Search Methods ^a

^a This is an extended version of the paper presented at the EWRL 2001 in Utrecht (The Netherlands). In this technical report, the derivation steps are presented with more detail, more footnotes, appendices and more (unfinished) ideas.

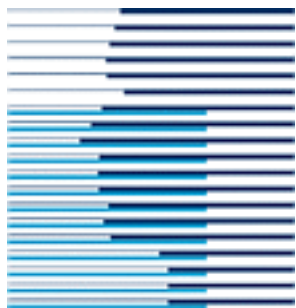
Ivo Kwee

Marcus Hutter

Jürgen Schmidhuber

IDSIA, Manno CH-6928, Switzerland.

{ivo,marcus,juergen}@idsia.ch



Technical Report No. IDSIA-14-01

November 2001

IDSIA / USI-SUPSI*

Instituto Dalle Molle di studi sull' intelligenza artificiale

Galleria 2

CH-6900 Manno, Switzerland

Gradient-based Reinforcement Planning in Policy-Search Methods *

Ivo Kwee Marcus Hutter Jürgen Schmidhuber
IDSIA, Manno CH-6928, Switzerland.
{ivo,marcus,juergen}@idsia.ch

November 2001

Abstract

We introduce a learning method called “gradient-based reinforcement planning” (GREP). Unlike traditional DP methods that improve their policy backwards in time, GREP is a gradient-based method that plans ahead and improves its policy *before* it actually acts in the environment. We derive formulas for the exact policy gradient that maximizes the expected future reward and confirm our ideas with numerical experiments.

1 Introduction

It has been shown that *planning* can dramatically improve convergence in reinforcement learning (RL) (?; ?). However, most RL methods that explicitly use planning that have been proposed are value (or Q -value) based methods, such as *Dyna-Q* or *prioritized sweeping*.

Recently, much attention is directed to so-called *policy-gradient* methods that improve their policy directly by calculating the derivative of the future expected reward with respect to the policy parameters. Gradient based methods are believed to be more advantageous than value-function based methods in huge state spaces and in POMDP settings. Probably the first gradient based RL formulation is class of REINFORCE algorithms of

Williams (?). Other more recent methods are, e.g., (?; ?; ?). Our approach of deriving the gradient has the flavor of (?) who derive the gradient using future state probabilities.

Our novel contribution in this paper is to combine gradient-based learning with explicit planning. We introduce “gradient-based reinforcement planning” (GREP) that improves a policy *before* actually interacting with the environment. We derive the formulas for the exact policy gradient and confirm our ideas in numerical experiments. GREP learns the action probabilities of a probabilistic policy for discrete problem. While we will illustrate GREP with a small MDP maze, it may be used for the hidden state in POMDPs.

*This is an extended version of the paper presented at the EWRL 2001 in Utrecht (The Netherlands). In this technical report, the derivation steps are presented with more detail, more footnotes, appendices and more (unfinished) ideas.

2 Derivation of the policy gradient

projection matrix

Let us denote the discrete space of *states* by $\mathcal{S} = \{1, \dots, N\}$. Our belief on \mathcal{S} is described by a probability vector \mathbf{s} of which each element s_i represents the probability of being in state i . We also define a set of actions $\mathcal{A} = 1, \dots, K$. The stochastic policy \mathcal{P} is represented by a matrix $\mathbf{P} : N \times K$ with elements $P_{ki} = p(a_k | s_i)$, i.e. the conditional probability of action k in state i .¹ Furthermore, let environment \mathcal{E} be defined by transition matrices \mathbf{T}_k ($k = 1, \dots, K$) with elements $T_{kji} = p(s_j | s_i, a_k)$, i.e. the transition probability to s_j in state s_i given action k .

Now we define the *projection matrix* \mathbf{F} with elements

$$F_{ji} = \sum_k T_{kji} P_{ki}. \tag{1}$$

Important is that matrix \mathbf{F} is *not* modelling the transition probabilities of the environment, but models the induced transition probability using policy \mathcal{P} in environment \mathcal{E} . The induced transition probability F_{ji} is a weighted sum over actions k of the transition probabilities T_{kij} with the policy parameters P_{ki} as the weights.

Expected state occupancy

Using the projection matrix \mathbf{F} , states \mathbf{s}_t and \mathbf{s}_{t+1} are related as $\mathbf{s}_{t+1} = \mathbf{F}\mathbf{s}_t$ and therefore $\mathbf{s}_t = \mathbf{F}^t \mathbf{s}_0$, where \mathbf{s}_0 is the state probability distribution at $t = 0$. We can now define the *expected state occupancy* as

$$\mathbf{z} = E[\mathbf{s} | \mathbf{s}_0] = \sum_{t=0}^{\infty} \gamma^t \mathbf{s}_t = \sum_{t=0}^{\infty} (\gamma \mathbf{F})^t \mathbf{s}_0 = (\mathbf{I} - \gamma \mathbf{F})^{-1} \mathbf{s}_0 \tag{2}$$

where γ is a discount factor in order to keep the sum finite. In the last step, we have recognized the sum as the Neumann representation of the inverse. Notice that \mathbf{z} is a solution of the linear equation

$$(\mathbf{I} - \gamma \mathbf{F}) \mathbf{z} = \mathbf{s}_0 \tag{3}$$

which is just the familiar Bellman equation for the expected occupancy probability \mathbf{z} .

Expected reward function

In *reinforcement learning* (RL) the objective is to maximize future reward. We define a reward vector \mathbf{r} in the same domain as \mathbf{s} . Using the expected occupancy \mathbf{z} the future expected reward H is simply

$$H = \langle \mathbf{r}, \mathbf{z} \rangle \tag{4}$$

where $\langle \cdot, \cdot \rangle$ is the scalar vector product.²

Because \mathbf{z} is a solution of Eq. 3 it is dependent on \mathbf{F} which in turn depends on policy \mathbf{P} . Given \mathbf{r} and \mathbf{s}_0 , our task is to find the optimal \mathbf{P}^* such that H is maximized, i.e.

$$\mathbf{P}^* = \arg \max_{\mathbf{P}} H. \tag{7}$$

¹For computational reasons, one often reparameterizes the policy using a Boltzmann distribution. Here in this paper the probability $p(a_k | s_i)$ is just given by P_{ki} and we do not use reparameterization in order to keep the analysis clear.

²In *optimal control* (OC) we want to reach some target state under some optimality conditions (mostly minimum time or minimum energy). We denote \mathbf{r} as our target distribution and denote the time-to-arrival as t^* . If t^* were known beforehand then we

$$H = \frac{1}{2} (\mathbf{r} - \mathbf{s}_{t^*})^2 = \frac{1}{2} (\mathbf{r} - \mathbf{F}^{t^*} \mathbf{s}_0)^2 \tag{5}$$

We can regard the calculation of the future expected reward as a composition of two operators

$$\mathcal{Q} : \mathbf{F} \mapsto \mathbf{z} \quad (8)$$

which maps the transition matrix \mathbf{F} to the expected occupancy probabilities \mathbf{z} , and

$$\mathcal{R} : \mathbf{z} \mapsto R \quad (9)$$

which maps the probabilities \mathbf{z} , given a reward distribution \mathbf{r} , to an expected reward value R .

Calculation of the policy gradient

A variation $\delta \mathbf{z}$ in the expected occupancy can be related to first order to a perturbation $\delta \mathbf{P}$ in the (stochastic) policy. To obtain the partial derivatives $\partial \mathbf{z} / \partial P_{ik}$, we differentiate Eq. 3 with respect to P_{ik} and obtain:

$$-\gamma \frac{\partial \mathbf{F}}{\partial P_{ik}} \mathbf{z} + (\mathbf{I} - \gamma \mathbf{F}) \frac{\partial \mathbf{z}}{\partial P_{ik}} = 0. \quad (10)$$

The right hand side of the equation is zero because we assume that \mathbf{s}_0 is independent of P_{ik} . Rearranging gives:

$$\frac{\partial \mathbf{z}}{\partial P_{ik}} = \gamma \mathbf{K} \frac{\partial \mathbf{F}}{\partial P_{ik}} \mathbf{z} \quad (11)$$

where $\mathbf{K} = (\mathbf{I} - \gamma \mathbf{F})^{-1}$.

From Eq. 4 and Eq. 10, together with the chain rule, we obtain the gradient of the RL error with respect to the policy parameters P_{ik} :

$$\frac{\partial H}{\partial P_{ik}} = \frac{\partial H}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial P_{ik}} = \left\langle \mathbf{r}, \gamma \mathbf{K} \frac{\partial \mathbf{F}}{\partial P_{ik}} \mathbf{z} \right\rangle = \gamma \left\langle \mathbf{K}^* \mathbf{r}, \frac{\partial \mathbf{F}}{\partial P_{ik}} \mathbf{z} \right\rangle \quad (12)$$

where A^* means the adjoint operator of A defined by $\langle u, Aw \rangle = \langle A^* u, w \rangle$. Let us define: $\mathbf{q} = \mathbf{K}^* \mathbf{r}$. While \mathbf{K} maps the initial state \mathbf{s}_0 to the future expected state occupancy \mathbf{z} , its adjoint, \mathbf{K}^* , maps the reward vector \mathbf{r} back to *expected reward* \mathbf{q} . The value of q_i represents the (pointwise) expected reward in state s_i for policy \mathbf{P} .³

Finally, differentiating Eq. 1 gives us $\partial \mathbf{F} / \partial P_{ik}$. Inserting this into Eq. 12 yields:

$$G_{ik} = \frac{\partial H}{\partial P_{ik}} \propto z_i \sum_j T_{kji} q_j. \quad (13)$$

In words, the gradient of H with respect to policy parameter P_{ik} (i.e. the probability of taking action a_k in state s_i) is proportional to the expected occupancy z_i times the weighted sum of expected reward q_j over next states ($j = 1, \dots, N$) weighted by the transition probabilities T_{kji} .

However the exact arrival time is mostly not known beforehand, the most we can do is to use minimize the (time-weighted) expected error

$$H = \sum_{t^*=0}^{\infty} \gamma^{t^*} \frac{1}{2} (\mathbf{r} - \mathbf{F}^{t^*} \mathbf{s}_0)^2 = \frac{1}{2} k \mathbf{r}^2 - \mathbf{r} \mathbf{z} + \frac{1}{2} \sum_{t^*=0}^{\infty} \gamma^{t^*} (\mathbf{F}^{t^*} \mathbf{s}_0)^2 \quad (6)$$

The first term on the right hand side is often not relevant because it is independent of \mathbf{F} . When we compare Eq. 6 with Eq. 4, we see that the OC error function has a quadratic term in \mathbf{F} which the RL error lacks.

³Indeed, this is a different way to define the traditional *value-function*. Note that generally, neither \mathbf{r} nor \mathbf{q} are probabilities because their 1-norm is generally not 1.

Note that the gradient could also have been approximated using finite differences which would need at least $1 + n^2$ field calculations.⁴ The adjoint method is much more efficient and needs only *two* field calculations.

Once we have the gradient \mathbf{G} , improving policy \mathbf{P} is now straight forward using gradient ascent or we can also use more sophisticated gradient-based methods such as nonlinear conjugate gradients (as in (?)). The optimization is nonlinear because \mathbf{z} and \mathbf{r} themselves depend on the current estimate of \mathbf{P} .

3 Computation of the optimal policy

We will introduce two algorithms that incorporate our ideas of gradient-based reinforcement planning. The first algorithm describes an off-line planning algorithm that finds the optimal policy but assumes that the environment transition probabilities are known. The second algorithm is an online version that could cope with unknown environments.

3.1 Offline GREP

If the environment transition probabilities T_{kji} are known, the agent may improve its policy using GREP. Our offline GREP planning algorithm consist of two steps:

1. *Plan ahead:* Compute the policy gradient \mathbf{G} in Eq. 12 and improve current policy

$$\mathbf{P} \leftarrow \mathbf{P} + \alpha \mathbf{G} \tag{14}$$

where α is a suitable step size parameter; for efficiency we can also perform a linesearch on α .

2. *Evaluate policy:* Repeat above until policy is optimal.

Matrix \mathbf{P} describes a probabilistic policy. We define the *maximum probable policy* (MPP) to be the deterministic policy by taking the maximum probable action at each state. It is not obvious that the MPP policy will converge to the global optimal solution but we expect MPP at least to be near-optimal.

Numerical experiments

We performed some numerical experiments using offline GREP. Our test problem was a pure planning task in a 10×10 toy maze (see Fig. 1) where the probabilistic policy \mathbf{P} represents the probability of taking a certain action at a certain maze position. The same figure also shows typical solutions for the quantities \mathbf{z} and \mathbf{q} , i.e. the expected occupancy and expected reward respectively (for certain \mathbf{P}).

After each GREP iteration, i.e. after each gradient calculation and \mathbf{P} update, we checked the obtained policy by running 20 simulations using the current value of \mathbf{P} . The probability weighted (PW) policy selects action k at state i proportional to P_{ik} , while the annealed PW policy uses an annealing factor of $T = 4$; we also simulated the MPP solution. Figure 2 shows the average simulated path length versus GREP iteration of the PW, the annealed PW policy and the derived MPP policy. In the left plot the initial policy \mathbf{P} was taken uniform. The right plot in the same

⁴ Finite difference approximation of the derivative $\partial \mathbf{z} / \partial T_{ij}$ involves computing \mathbf{z} for \mathbf{F} and then perturbing a single T_{ij} in \mathbf{F} by a tiny amount dT and subsequently recomputing \mathbf{z}' . Then the derivative is approximated by $\partial \mathbf{z} / \partial T_{ij} \approx (\mathbf{z}' - \mathbf{z}) / dT$. For a $n \times n$ matrix \mathbf{F} , one would need to repeat this for every element and would require a total upto $1 + n^2$ calculations of \mathbf{z} .

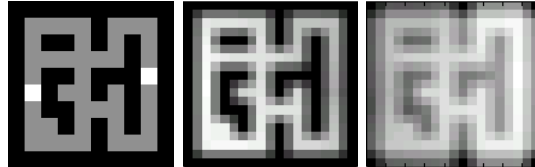


Figure 1: Left: 10×10 toy maze with start at left and goal at right side. Center: plot of expected occupancy \mathbf{z} . Right: plot of expected reward \mathbf{q} . White corresponds to higher probability. [Blurring is due to visualisation only].

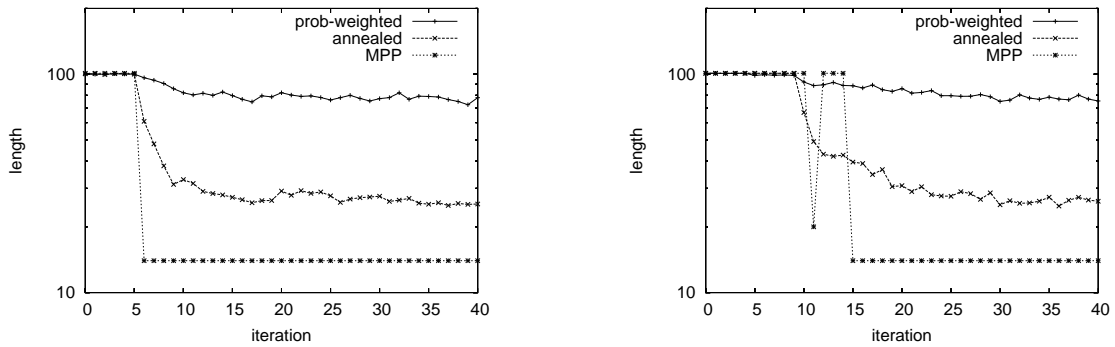


Figure 2: Plot of simulated path length versus GREP iteration of a small toy MDP maze for the probability weighted (PW) policy, annealed PW policy and MPP policy. The shortest path to goal is 14. Left: starting from initial uniform policy. Right: starting from initial random policy.

figure shows the simulated path lengths from a random policy; also here the MPP finds the optimal solution but slightly later.

We see from the figure that in both cases the probability-weighted (PW) policy is improving during the GREP iterations. However, the convergence is very slow which shows the severe non-linearity of the problem. The annealed PW does perform better than PW. Finally, we see that MPP finds the optimal solution quickly within a few iterations. Using Dijkstra’s method, we confirmed that the found MPP policy was in agreement with the global shortest path solution.

Online GREP

The account below describe an idea to use GREP when the environment is not known beforehand⁵. The steps actually interleave “Kalman filter”-like estimation of the unknown environment transition probabilities with the explicit planning of GREP. In fact, it also includes a step to estimate a possibly unknown (linear) sensor mapping. Apart from the policy matrix \mathbf{P} , we need to estimate also the (environment) transition probabilities T_{kji} and possibly sensor matrix \mathbf{B} . We can optimize for all parameters by iteratively ascending to their conditional mode. The conditional maximizing steps are easy:

1. *Plan ahead:* Compute the policy gradient \mathbf{G} in Eq. 12 and improve current policy

$$\mathbf{P} \leftarrow \mathbf{P} + \alpha \mathbf{G} \tag{15}$$

where α is a suitable step size parameter; for efficiency we can also perform a linesearch. After, we need to renormalize the columns of \mathbf{P} . See note below on policy regularization.

2. *Select action:* Given state estimate \mathbf{s}_t , draw an action k from the policy according to:

$$k \sim \mathbf{P}_t \mathbf{s}_t.$$

and receive reward R and estimate new state \mathbf{s}_{t+1} .

3. *Estimate state:* Observe \mathbf{y}_{t+1} and estimate \mathbf{s}_{t+1} using

$$p(\mathbf{s}_{t+1} | \mathbf{y}_{t+1}, \mathbf{s}_t, k) \propto p(\mathbf{y}_{t+1} | \mathbf{s}_{t+1}) p(\mathbf{s}_{t+1} | \mathbf{F}_k, \mathbf{s}_t). \tag{16}$$

Assuming Gaussian noise for observations and state estimates we obtain:

$$\mathbf{s}_{t+1} = \frac{\mathbf{H}_s \mathbf{F}_k \mathbf{s}_t + \mathbf{B}^{-1} \mathbf{H}_y \mathbf{y}_{t+1}}{\mathbf{H}_s + \mathbf{B}^{-1} \mathbf{H}_y \mathbf{B}^{-T}} \tag{17}$$

where \mathbf{B} is the *sensor matrix* that maps internal state \mathbf{s} to observations \mathbf{y} . Matrices \mathbf{H}_s and \mathbf{H}_y are the inverse covariance (or so called *precisions*) of state \mathbf{s} and observation \mathbf{y} respectively.

4. *Estimate sensors:* In case also sensor matrix \mathbf{B} is unknown, we have to perform an additional estimation step for \mathbf{B} . This is common step in the standard Kalman formulation.
5. *Estimate environment:* Given action k and reward R , we update the reward vector

$$r_j \leftarrow R \tag{18}$$

and the environment transition probabilities

$$d\mathbf{F}_k \propto (\mathbf{s}_{t+1} - \mathbf{F}_k \mathbf{s}_t) \mathbf{s}_t^T (\mathbf{s}_t \mathbf{s}_t^T)^{-1}, \tag{19}$$

⁵At the time of writing, we have not implemented this idea yet

or if $(\mathbf{s}_t \mathbf{s}_t^T)^{-1}$ does not exist we can use

$$d\mathbf{F}_k \propto (\mathbf{s}_{t+1} - \mathbf{F}_k \mathbf{s}_t)(\mathbf{s}_t^T \mathbf{s}_t)^{-1} \mathbf{s}_t^T. \quad (20)$$

where $\mathbf{s}_t = j$, and reestimate the environment transition probabilities

$$\mathbf{T}_k \leftarrow \mathbf{T}_k + (\mathbf{s}_{t+1} - \mathbf{T}_k \mathbf{s}_t) \mathbf{s}_t^T. \quad (21)$$

After the update, one should set entries in \mathbf{F}'_k that corresponds to physically impossible transitions to zero. After, we need to renormalize the columns of \mathbf{T}_k . It is important to note that given \mathbf{s}_t and \mathbf{s}_{t+1} transition matrix \mathbf{T}_k is conditionally independent of the policy \mathbf{P} . That is to say, we can obtain an accurate model of the environment using, e.g., just a random walk.

6. Repeat 1.

To draw a picture of what is happening. In the planning stage, based on the current (and maybe not accurate) environment model, the agent tries to improve its current policy by planning ahead using the gradient in Eq. 13. Remember that the gradient involves simulating paths from the current state and adjoint paths from the goal. In the action stage the agent samples an action from its policy. Then the agent senses the new state and updates its environment model using this new information. Notice that policy improvement is not done “backwards” as traditionally is done in DP methods but “forward” by planning ahead.

4 Conclusions

Future topics

We have tacitly assumed that \mathbf{z} and \mathbf{q} are computed using the same discount factor γ . However, we could introduce separate parameters γ_z and γ_q which effectively assigns a different “forward time window” for \mathbf{z} and a “backward time window” for \mathbf{q} . In fact when $\gamma_z \rightarrow 0$ we have a “one-step-look-ahead”. Alternatively, in the limit of $\gamma_q \rightarrow 0$ we obtain a gradient for a greedy policy that maximize only “immediate reward”. How both parameters affect GREP’s performance is a topic for future research.

The above suggests that GREP can be viewed as a generalization to “one-step-look-ahead” policy improvement. In fact, a “one-step-look-ahead” improvement rule using can be obtained for $\gamma_z \rightarrow 0$ simply by taking $\mathbf{z} = \mathbf{s}_t$ in Eq. 13. Such an approach would be “policy greedy” in a sense that it updates the policy only locally. We expect GREP to perform better because it updates the policy more globally; whether this in fact improves GREP is also a remaining issue for future research.

The interleaving of GREP with a Kalman-like estimation procedure of the environment could handle a variety of interesting problems such as planning in POMDP environments.

We must mention that appropriate reparameterization of the stochastic policy, e.g. using a Boltzman distribution, could improve the convergence. We have not pursued this further.

Summary

We have introduced a learning method called “gradient-based reinforcement planning” (GREP). GREP needs a model of the environment and plans ahead to improve its policy *before* it actually acts in the environment. We have derived formulas for the exact policy gradient.

Numerical experiments suggest that the probabilistic policy indeed converges to an optimal policy—but quite slowly. We found that (at least in our toy example) the optimal solution can be found much faster by annealing or simply by taking the most probable action at each state.

Further work will be to incorporate GREP in online RL learning tasks where the environment parameters, i.e. transition probabilities T_{kji} , are unknown and have to be learned. While an analytical solution for \mathbf{q} and \mathbf{z} are only viable for small problem sizes, for larger problems we probably need to investigate Monte Carlo or DP methods.

References

- Baird][1998]baird98gradient Baird, L. C. (1998). Gradient descent for general reinforcement learning. *Advances in Neural Information Processing Systems*. MIT Press.
- Baxter and Bartlett][1999]baxter99direct Baxter, J., & Bartlett, P. (1999). *Direct gradient-based reinforcement learning: I. gradient estimation algorithms* (Technical Report). Research School of Information Sciences and Engineering, Australian National University.
- Difilippo et al.][1996]difilippo Difilippo, F. C., Goldstein, M., Worley, B. A., & Ryman, J. C. (1996). Adjoint Monte Carlo methods for radiotherapy treatment planning. *Trans. Am. Nucl. Soc.*, 74, 14–16.
- Ng et al.][1999]ng99policy Ng, A., Parr, R., & Koller, D. (1999). Policy search via density estimation. *Advances in Neural Information Processing Systems*. MIT Press.
- Schmidhuber][1990]Schmidhuber:90sandiego Schmidhuber, J. (1990). An on-line algorithm for dynamic reinforcement learning and planning in reactive environments. *Proc. IEEE/INNS International Joint Conference on Neural Networks, San Diego* (pp. 253–258).
- Sutton and Barto][1998]suttonbarto Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning. An introduction*. MIT Press, Cambridge.
- Sutton et al.][2000]sutton00policy Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. *Advances in Neural Information Processing Systems*. MIT Press.
- Williams][1992]williams92simple Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8, 229–256.

Appendix A: Implicit policies

In deterministic environments where the state-action pair (s_i, a_m) uniquely leads to a state s_j , i.e. $T_{kji} = \delta_{km}$ ($k = 1, \dots, K$) the projection \mathbf{F} is solely determined by the policy \mathbf{z} , and *vice versa*. We refer to this as the case of *implicit policy* because the policy is implicitly implied in the induced transition probability T_{ji} .

In such environments we can suffice to solve for \mathbf{F} directly and omit parameterization through \mathbf{z} . From Eq. 1 we see that

$$P_{im} = T_{ji} \tag{22}$$

and using a similar derivation as we have done for $\partial H / \partial P_{ik}$, it can be shown that the gradient of H with respect to \mathbf{F} is given by

$$\mathbf{G} = \mathbf{r}\mathbf{z}^T. \tag{23}$$

An important point must be mentioned. In most cases many elements T_{ji} are zero, representing an absent transition between S_i and S_j . Naively updating \mathbf{F} using the full gradient \mathbf{G} would incur complete fill-in of \mathbf{F} which is in most cases not desirable or even physically incorrect. Therefore, one must check the gradient each time and set impossible transition probabilities to zero. We will refer to this “heuristically corrected” gradient as $\tilde{\mathbf{G}}$. Also, after each update, we have to renormalize the columns of \mathbf{F} . The rank-one update in Eq. 23 is interesting because it provides an efficient means of calculating the inverse in Eq. 2.

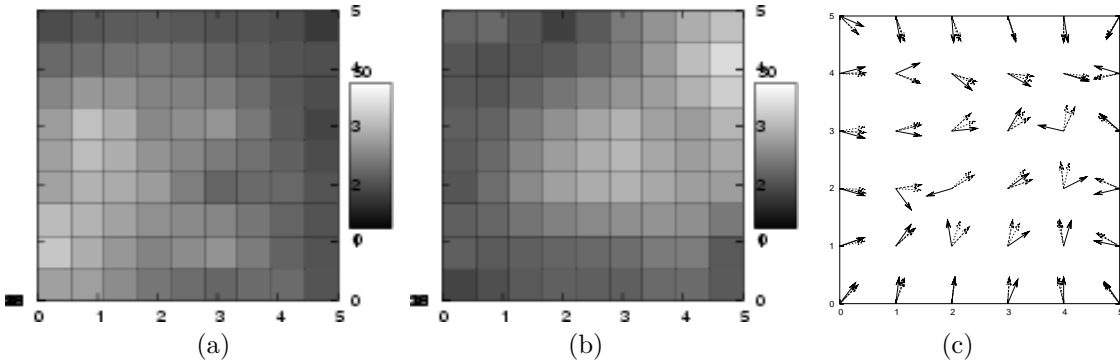


Figure 3: MC calculations of a 6×6 toy maze. The agent is at $(0,1)$ and targets $(5,4)$. (a) Expected occupancy, (b) adjoint probability, and (c) normalized policy gradient for $n = 4, n = 40, n = 200$. Each vector is computed as $\mathbf{v} = \sum_k P_{ik} \mathbf{e}_k$ where \mathbf{e}_k is the unit vector along the state change induced by action a_k .

Appendix B: Monte Carlo gradient sampling

In our example, we calculated \mathbf{z} and \mathbf{q} in Eq.2 by linear programming. For large state spaces the matrix inversion quickly becomes too computationally intensive and probably traditional dynamic programming based methods would be more efficient.

Instead, we investigated to use Monte Carlo (MC) simulation. We use *forward sampling* to approximate the expected state occupancies in \mathbf{z} and use, so-called, *adjoint Monte Carlo sampling* (?) to estimate the adjoint reward \mathbf{q} . Adjoint MC simulation is far more efficient than would we have estimated each q_i by a separate MC run.⁶ By performing the simulation backward from \mathbf{r} , we obtain all values of \mathbf{q} using only a single MC run.

Fig. 3 shows the MC approximations of \mathbf{z} and \mathbf{q} . On the right of the same figure, we have plotted the computed policy-gradient based on MC estimates using a minimum number of $n = \{20, 40, 200\}$ samples. To compare them with the exact gradient, we calculated the exact values of \mathbf{z} and \mathbf{q} by inverting the linear system. For larger number of samples, the gradient vector do indeed point more strongly towards the goal.

An important feature of general Monte Carlo methods is that they automatically concentrate their sampling to the important regions of the parameter space —mostly proportional to the posterior or the likelihood. For our purpose of sampling the gradient, to even more concentrate the sampling density towards the regions of large gradient values, we have tried to apply *annealing*. To sample from a density $p(\theta)$ we may sample from the annealed function $p_\gamma(\theta) = p(\theta)^\gamma / \int p(\theta)^\gamma d\theta$ and reweight each sample with its importance weight $1/p_\gamma(\theta)$. For $\gamma \rightarrow \infty$, the set of samples converges to the maximum probable gradient.

In conclusion, our approach of separately estimating \mathbf{q} and \mathbf{z} using MC and *then* (elementwise) multiply their solutions, doesn't really brought clear advantages. If we could sample from the joint distribution $q_i z_i$ (i.e. elementwise product) then MC would clearly turn out to be a very efficient method.

⁶With one MC run we mean performing, e.g., 10000 trials from a fixed state.