

IMPROVING POLICY GRADIENT BY EXPLORING UNDER-APPRECIATED REWARDS

Ofir Nachum*, Mohammad Norouzi, Dale Schuurmans†
 Google Brain
 {ofirnachum, mnorouzi, schuurmans}@google.com

ABSTRACT

This paper presents a novel form of policy gradient for model-free reinforcement learning (RL) with improved exploration properties. Current policy-based methods use entropy regularization to encourage undirected exploration of the reward landscape, which is ineffective in high dimensional spaces with sparse rewards. We propose a more directed exploration strategy that promotes exploration of *under-appreciated reward* regions. An action sequence is considered under-appreciated if its log-probability under the current policy under-estimates its resulting reward. The proposed exploration strategy is easy to implement, requiring small modifications to the REINFORCE algorithm. We evaluate the approach on a set of algorithmic tasks that have long challenged RL methods. Our approach reduces hyper-parameter sensitivity and demonstrates significant improvements over baseline methods. The proposed algorithm successfully solves a benchmark multi-digit addition task and generalizes to long sequences, which, to our knowledge, is the first time that a pure RL method has solved addition using only reward feedback.

1 INTRODUCTION

Humans can reason about symbolic objects and solve algorithmic problems. After learning to count and then manipulate numbers via simple arithmetic, people eventually learn to invent new algorithms and even reason about their correctness and efficiency. The ability to invent new algorithms is fundamental to artificial intelligence (AI). Although symbolic reasoning has a long history in AI (Russell et al., 2003), only recently have statistical machine learning and neural network approaches begun to make headway in automated algorithm discovery (Reed & de Freitas, 2016; Kaiser & Sutskever, 2016; Neelakantan et al., 2016), which would constitute an important milestone on the path to AI. Nevertheless, most of the recent successes depend on the use of strong supervision to learn a mapping from a set of training inputs to outputs by maximizing a conditional log-likelihood, very much like neural machine translation systems (Sutskever et al., 2014; Bahdanau et al., 2015). Such a dependence on strong supervision is a significant limitation that does not match the ability of people to invent new algorithmic procedures based solely on trial and error.

By contrast, *reinforcement learning (RL)* methods (Sutton & Barto, 1998) hold the promise of searching over discrete objects such as symbolic representations of algorithms by considering much weaker feedback in the form of a simple verifier that tests the correctness of a program execution on a given problem instance. Despite the recent excitement around the use of RL to tackle Atari games (Mnih et al., 2015) and Go (Silver et al., 2016), standard RL methods are not yet able to consistently and reliably solve algorithmic tasks in all but the simplest cases (Zaremba & Sutskever, 2014). A key property of algorithmic problems that makes them challenging for RL is *reward sparsity*, *i.e.*, a policy usually has to get a long action sequence exactly right to obtain a non-zero reward.

*Work done as a member of the Google Brain Residency program (g.co/brainresidency)

†Also at the Department of Computing Science, University of Alberta, daes@ualberta.ca

We believe one of the key factors limiting the effectiveness of current RL methods in a sparse reward setting is the use of *undirected exploration* strategies (Thrun, 1992), such as ϵ -greedy and entropy regularization (Williams & Peng, 1991). For long action sequences with delayed sparse reward, it is hopeless to explore the space uniformly and blindly. Instead, we propose a formulation to encourage exploration of action sequences that are *under-appreciated* by the current policy. Our formulation considers an action sequence to be under-appreciated if the model’s log-probability assigned to an action sequence under-estimates the resulting reward from the action sequence. Exploring under-appreciated states and actions encourages the policy to have a better calibration between its log-probabilities and observed reward values, even for action sequences with negligible rewards. This effectively increases exploration around neglected action sequences.

We term our proposed technique *under-appreciated reward exploration (UREX)*. We show that the objective given by UREX is a combination of a mode seeking objective (standard REINFORCE) and a mean seeking term, which provides a well motivated trade-off between exploitation and exploration. To empirically evaluate our method, we take a set of algorithmic tasks such as sequence reversal, multi-digit addition, and binary search. We choose to focus on these tasks because, although simple, they present a difficult sparse reward setting which has limited the success of standard RL approaches. The experiments demonstrate that UREX significantly outperforms baseline RL methods, such as entropy regularized REINFORCE and one-step Q-learning, especially on the more difficult tasks, such as multi-digit addition. Moreover, UREX is shown to be more robust to changes of hyper-parameters, which makes hyper-parameter tuning less tedious in practice. In addition to introducing a new variant of policy gradient with improved performance, our paper is the first to demonstrate strong results for an RL method on algorithmic tasks. To our knowledge, the addition task has not been solved by any model-free reinforcement learning approach. We observe that some of the policies learned by UREX can successfully generalize to long sequences; *e.g.*, in 2 out of 5 random restarts, the policy learned by UREX for the addition task correctly generalizes to addition of numbers with 2000 digits with no mistakes, even though training sequences are at most 33 digits long.

2 NEURAL NETWORKS FOR LEARNING ALGORITHMS

Although research on using neural networks to learn algorithms has witnessed a surge of recent interest, the problem of program induction from examples has a long history in many fields, including program induction, inductive logic programming (Lavrac & Dzeroski, 1994), relational learning (Kemp et al., 2007) and regular language learning (Angulin, 1987). Rather than presenting a comprehensive survey of program induction here, we focus on neural network approaches to algorithmic tasks and highlight the relative simplicity of our neural network architecture.

Most successful applications of neural networks to algorithmic tasks rely on strong supervision, where the inputs and target outputs are completely known *a priori*. Given a dataset of examples, one learns the network parameters by maximizing the conditional likelihood of the outputs via back-propagation (*e.g.*, Reed & de Freitas (2016); Kaiser & Sutskever (2016); Vinyals et al. (2015)). However, target outputs may not be available for novel tasks, for which no prior algorithm is known to be available. A more desirable approach to inducing algorithms, followed in this paper, advocates using self-driven learning strategies that only receive reinforcement based on the outputs produced. Hence, just by having access to a verifier for an algorithmic problem, one can aim to learn an algorithm. For example, if one does not know how to sort an array, but can check the extent to which an array is sorted, then one can provide the reward signal necessary for learning sorting algorithms.

We formulate learning algorithms as an RL problem and make use of model-free policy gradient methods to optimize a set parameters associated with the algorithm. In this setting, the goal is to learn a policy π_θ that given an observed state \mathbf{s}_t at step t , estimates a distribution over the next action \mathbf{a}_t , denoted $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$. Actions represent the commands within the algorithm and states represent the joint state of the algorithm and the environment. Previous work in this area has focused on augmenting a neural network with additional structure and increased capabilities (Zaremba & Sutskever, 2015; Graves et al., 2016). In contrast, we utilize a simple architecture based on a standard recurrent neural network (RNN) with LSTM cells (Hochreiter & Schmidhuber, 1997) as depicted in Figure 1. At each episode, the environment is initialized with a latent state \mathbf{h} , unknown to the agent, which determines \mathbf{s}_1 and the subsequent state transition and reward functions. Once the agent observes \mathbf{s}_1

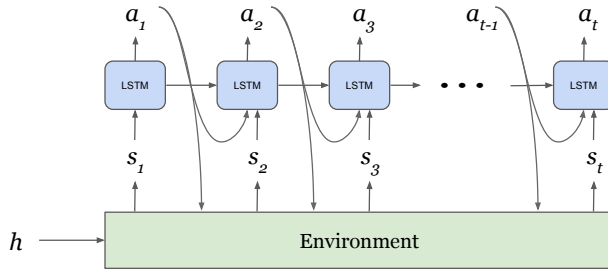


Figure 1: The agent’s RNN architecture that represents a policy. The environment is initialized with a latent vector \mathbf{h} . At time step t , the environment produces a state \mathbf{s}_t , and the agent takes as input \mathbf{s}_t and the previously sampled action \mathbf{a}_{t-1} and produces a distribution over the next action $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$. Then, we sample a new action \mathbf{a}_t and apply it to the environment.

as the input to the RNN, the network outputs a distribution $\pi_\theta(\mathbf{a}_1 | \mathbf{s}_1)$, from which an action \mathbf{a}_1 is sampled. This action is applied to the environment, and the agent receives a new state observation \mathbf{s}_2 . The state \mathbf{s}_2 and the previous action \mathbf{a}_1 are then fed into the RNN and the process repeats until the end of the episode. Upon termination, a reward signal is received.

3 LEARNING A POLICY BY MAXIMIZING EXPECTED REWARD

We start by discussing the most common form of policy gradient, REINFORCE (Williams, 1992), and its entropy regularized variant (Williams & Peng, 1991). REINFORCE has been applied to model-free policy-based learning with neural networks and algorithmic domains (Zaremba & Sutskever, 2015; Graves et al., 2016).

The goal is to learn a policy π_θ that, given an observed state \mathbf{s}_t at step t , estimates a distribution over the next action \mathbf{a}_t , denoted $\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)$. The environment is initialized with a latent vector, \mathbf{h} , which determines the initial observed state $\mathbf{s}_1 = g(\mathbf{h})$, and the transition function $\mathbf{s}_{t+1} = f(\mathbf{s}_t, \mathbf{a}_t | \mathbf{h})$. Note that the use of nondeterministic transitions f as in Markov decision processes (MDP) may be recovered by assuming that \mathbf{h} includes the random seed for the any nondeterministic functions. Given a latent state \mathbf{h} , and $\mathbf{s}_{1:T} \equiv (\mathbf{s}_1, \dots, \mathbf{s}_T)$, the model probability of an action sequence $\mathbf{a}_{1:T} \equiv (\mathbf{a}_1, \dots, \mathbf{a}_T)$ is expressed as,

$$\pi_\theta(\mathbf{a}_{1:T} | \mathbf{h}) = \prod_{t=1}^T \pi_\theta(\mathbf{a}_t | \mathbf{s}_t), \quad \text{where } \mathbf{s}_1 = g(\mathbf{h}), \mathbf{s}_{t+1} = f(\mathbf{s}_t, \mathbf{a}_t | \mathbf{h}) \text{ for } 1 \leq t < T.$$

The environment provides a reward at the end of the episode, denoted $r(\mathbf{a}_{1:T} | \mathbf{h})$. For ease of readability we drop the subscript from $\mathbf{a}_{1:T}$ and simply write $\pi_\theta(\mathbf{a} | \mathbf{h})$ and $r(\mathbf{a} | \mathbf{h})$.

The objective used to optimize the policy parameters, θ , consists of maximizing expected reward under actions drawn from the policy, plus an optional maximum entropy regularizer. Given a distribution over initial latent environment states $p(\mathbf{h})$, we express the regularized expected reward as,

$$\mathcal{O}_{\text{RL}}(\theta; \tau) = \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h})} \left\{ \sum_{\mathbf{a} \in \mathcal{A}} \pi_\theta(\mathbf{a} | \mathbf{h}) \left[r(\mathbf{a} | \mathbf{h}) - \tau \log \pi_\theta(\mathbf{a} | \mathbf{h}) \right] \right\}. \quad (1)$$

When π_θ is a non-linear function defined by a neural network, finding the global optimum of θ is challenging, and one often resorts to gradient-based methods to find a local optimum of $\mathcal{O}_{\text{RL}}(\theta; \tau)$. Given that $\frac{d}{d\theta} \pi_\theta(\mathbf{a}) = \pi_\theta(\mathbf{a}) \frac{d}{d\theta} \log \pi_\theta(\mathbf{a})$ for any \mathbf{a} such that $\pi_\theta(\mathbf{a}) > 0$, one can verify that,

$$\frac{d}{d\theta} \mathcal{O}_{\text{RL}}(\theta; \tau | \mathbf{h}) = \sum_{\mathbf{a} \in \mathcal{A}} \pi_\theta(\mathbf{a} | \mathbf{h}) \frac{d}{d\theta} \log \pi_\theta(\mathbf{a} | \mathbf{h}) \left[r(\mathbf{a} | \mathbf{h}) - \tau \log \pi_\theta(\mathbf{a} | \mathbf{h}) - \tau \right]. \quad (2)$$

Because the space of possible actions \mathcal{A} is large, enumerating over all of the actions to compute this gradient is infeasible. Williams (1992) proposed to compute the stochastic gradient of the expected

reward by using Monte Carlo samples. Using Monte Carlo samples, one first draws N *i.i.d.* samples from the latent environment states $\{\mathbf{h}^{(n)}\}_{n=1}^N$, and then draws K *i.i.d.* samples $\{\mathbf{a}^{(k)}\}_{k=1}^K$ from $\pi_\theta(\mathbf{a} | \mathbf{h}^{(n)})$ to approximate the gradient of (1) by using (2) as,

$$\frac{d}{d\boldsymbol{\theta}} \mathcal{O}_{\text{RL}}(\boldsymbol{\theta}; \tau) \approx \frac{1}{NK} \sum_{n=1}^N \sum_{k=1}^K \frac{d}{d\boldsymbol{\theta}} \log \pi_\theta(\mathbf{a}^{(k)} | \mathbf{h}^{(n)}) \left[r(\mathbf{a}^{(k)} | \mathbf{h}^{(n)}) - \tau \log \pi_\theta(\mathbf{a}^{(k)} | \mathbf{h}^{(n)}) - \tau \right]. \quad (3)$$

This reparametrization of the gradients is the key to the REINFORCE algorithm. To reduce the variance of (3), one uses rewards \hat{r} that are shifted by some offset values,

$$\hat{r}(\mathbf{a}^{(k)} | \mathbf{h}) = r(\mathbf{a}^{(k)} | \mathbf{h}) - b(\mathbf{h}), \quad (4)$$

where b is known as a *baseline* or sometimes called a *critic*. Note that subtracting any offset from the rewards in (1) simply results in shifting the objective \mathcal{O}_{RL} by a constant.

Unfortunately, directly maximizing expected reward (*i.e.*, when $\tau = 0$) is prone to getting trapped in a local optimum. To combat this tendency, Williams & Peng (1991) augmented the expected reward objective by including a maximum entropy regularizer ($\tau > 0$) to promote greater exploration. We will refer to this variant of REINFORCE as MENT (maximum entropy exploration).

4 UNDER-APPRECIATED REWARD EXPLORATION (UREX)

To explain our novel form of policy gradient, we first note that the optimal policy π_τ^* , which globally maximizes $\mathcal{O}_{\text{RL}}(\boldsymbol{\theta}; \tau | \mathbf{h})$ in (1) for any $\tau > 0$, can be expressed as,

$$\pi_\tau^*(\mathbf{a} | \mathbf{h}) = \frac{1}{Z(\mathbf{h})} \exp \left\{ \frac{1}{\tau} r(\mathbf{a} | \mathbf{h}) \right\}, \quad (5)$$

where $Z(\mathbf{h})$ is a normalization constant making π_τ^* a distribution over the space of action sequences \mathcal{A} . One can verify this by first acknowledging that,

$$\mathcal{O}_{\text{RL}}(\boldsymbol{\theta}; \tau | \mathbf{h}) = -\tau D_{\text{KL}}(\pi_\theta(\cdot | \mathbf{h}) \| \pi_\tau^*(\cdot | \mathbf{h})). \quad (6)$$

Since $D_{\text{KL}}(p \| q)$ is non-negative and zero iff $p = q$, then π_τ^* defined in (5) maximizes \mathcal{O}_{RL} . That said, given a particular form of π_θ , finding θ that exactly characterizes π_τ^* may not be feasible.

The KL divergence $D_{\text{KL}}(\pi_\theta \| \pi_\tau^*)$ is known to be mode seeking (Murphy, 2012, Section 21.2.2) even with entropy regularization ($\tau > 0$). Learning a policy by optimizing this direction of the KL is prone to falling into a local optimum resulting in a sub-optimal policy that omits some of the modes of π_τ^* . Although entropy regularization helps mitigate the issues as confirmed in our experiments, it is not an effective exploration strategy as it is undirected and requires a small regularization coefficient τ to avoid too much random exploration. Instead, we propose a directed exploration strategy that improves the mean seeking behavior of policy gradient in a principled way.

We start by considering the alternate mean seeking direction of the KL divergence, $D_{\text{KL}}(\pi_\tau^* \| \pi_\theta)$. Norouzi et al. (2016) considered this direction of the KL to directly learn a policy by optimizing

$$\mathcal{O}_{\text{RAML}}(\boldsymbol{\theta}; \tau) = \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h})} \left\{ \tau \sum_{\mathbf{a} \in \mathcal{A}} \pi_\tau^*(\mathbf{a} | \mathbf{h}) \log \pi_\theta(\mathbf{a} | \mathbf{h}) \right\}, \quad (7)$$

for structured prediction. This objective has the same optimal solution π_τ^* as \mathcal{O}_{RL} since,

$$\mathcal{O}_{\text{RAML}}(\boldsymbol{\theta}; \tau | \mathbf{h}) = -\tau D_{\text{KL}}(\pi_\tau^*(\cdot | \mathbf{h}) \| \pi_\theta(\cdot | \mathbf{h})) + \text{const}. \quad (8)$$

Norouzi et al. (2016) argue that in some structured prediction problems when one can draw samples from π_τ^* , optimizing (7) is more effective than (1), since no sampling from a non-stationary policy π_θ is required. If π_θ is a log-linear model of a set of features, $\mathcal{O}_{\text{RAML}}$ is convex in $\boldsymbol{\theta}$ whereas \mathcal{O}_{RL} is not, even in the log-linear case. Unfortunately, in scenarios that the reward landscape is unknown or computing the normalization constant $Z(\mathbf{h})$ is intractable, sampling from π_τ^* is not straightforward.

In RL problems, the reward landscape is completely unknown, hence sampling from π_τ^* is intractable. This paper proposes to approximate the expectation with respect to π_τ^* by using *self-normalized importance sampling* (Owen, 2013), where the proposal distribution is π_θ and the reference distribution is π_τ^* . For importance sampling, one draws K *i.i.d.* samples $\{\mathbf{a}^{(k)}\}_{k=1}^K$ from

$\pi_\theta(\mathbf{a} \mid \mathbf{h})$ and computes a set of normalized importance weights to approximate $\mathcal{O}_{\text{RAML}}(\boldsymbol{\theta}; \tau \mid \mathbf{h})$ as,

$$\mathcal{O}_{\text{RAML}}(\boldsymbol{\theta}; \tau \mid \mathbf{h}) \approx \tau \sum_{k=1}^K \frac{w_\tau(\mathbf{a}^{(k)} \mid \mathbf{h})}{\sum_{m=1}^K w_\tau(\mathbf{a}^{(m)} \mid \mathbf{h})} \log \pi_\theta(\mathbf{a}^{(k)} \mid \mathbf{h}), \quad (9)$$

where $w_\tau(\mathbf{a}^{(k)} \mid \mathbf{h}) \propto \pi_\tau^*/\pi_\theta$ denotes an importance weight defined by,

$$w_\tau(\mathbf{a}^{(k)} \mid \mathbf{h}) = \exp \left\{ \frac{1}{\tau} r(\mathbf{a}^{(k)} \mid \mathbf{h}) - \log \pi_\theta(\mathbf{a}^{(k)} \mid \mathbf{h}) \right\}. \quad (10)$$

One can view these importance weights as evaluating the discrepancy between scaled rewards r/τ and the policy’s log-probabilities $\log \pi_\theta$. Among the K samples, a sample that is least appreciated by the model, *i.e.*, has the largest $r/\tau - \log \pi_\theta$, receives the largest positive feedback in (9).

In practice, we have found that just using the importance sampling RAML objective in (9) does not always yield promising solutions. Particularly, at the beginning of training, when π_θ is still far away from π_τ^* , the variance of importance weights is too large, and the self-normalized importance sampling procedure results in poor approximations. To stabilize early phases of training and ensure that the model distribution π_θ achieves large expected reward scores, we combine the expected reward and RAML objectives to benefit from the best of their mode and mean seeking behaviors. Accordingly, we propose the following objective that we call *under-appreciated reward exploration (UREX)*,

$$\mathcal{O}_{\text{UREX}}(\boldsymbol{\theta}; \tau) = \mathbb{E}_{\mathbf{h} \sim p(\mathbf{h})} \left\{ \sum_{\mathbf{a} \in \mathcal{A}} \left[\pi_\theta(\mathbf{a} \mid \mathbf{h}) r(\mathbf{a} \mid \mathbf{h}) + \tau \pi_\tau^*(\mathbf{a} \mid \mathbf{h}) \log \pi_\theta(\mathbf{a} \mid \mathbf{h}) \right] \right\}, \quad (11)$$

which is the sum of the expected reward and RAML objectives. In our preliminary experiments, we considered a composite objective of $\mathcal{O}_{\text{RL}} + \mathcal{O}_{\text{RAML}}$, but we found that removing the entropy term is beneficial. Hence, the $\mathcal{O}_{\text{UREX}}$ objective does not include entropy regularization. Accordingly, the optimum policy for $\mathcal{O}_{\text{UREX}}$ is no longer π_τ^* , as it was for \mathcal{O}_{RL} and $\mathcal{O}_{\text{RAML}}$. Appendix A derives the optimal policy for $\mathcal{O}_{\text{UREX}}$ as a function of the optimal policy for \mathcal{O}_{RL} . We find that the optimal policy of UREX is more sharply concentrated on the high reward regions of the action space, which may be an advantage for UREX, but we leave more analysis of this behavior to future work.

To compute the gradient of $\mathcal{O}_{\text{UREX}}(\boldsymbol{\theta}; \tau)$, we use the self-normalized importance sampling estimate outlined in (9). We assume that the importance weights are constant and contribute no gradient to $\frac{d}{d\boldsymbol{\theta}} \mathcal{O}_{\text{UREX}}(\boldsymbol{\theta}; \tau)$. To approximate the gradient, one draws N *i.i.d.* samples from the latent environment states $\{\mathbf{h}^{(n)}\}_{n=1}^N$, and then draws K *i.i.d.* samples $\{\mathbf{a}^{(k)}\}_{k=1}^K$ from $\pi_\theta(\mathbf{a} \mid \mathbf{h}^{(n)})$ to obtain

$$\frac{d}{d\boldsymbol{\theta}} \mathcal{O}_{\text{UREX}}(\boldsymbol{\theta}; \tau) \approx \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \frac{d}{d\boldsymbol{\theta}} \log \pi_\theta(\mathbf{a}^{(k)} \mid \mathbf{h}^{(n)}) \left[\frac{1}{K} \hat{r}(\mathbf{a}^{(k)} \mid \mathbf{h}^{(n)}) + \tau \frac{w_\tau(\mathbf{a}^{(k)} \mid \mathbf{h}^{(n)})}{\sum_{m=1}^K w_\tau(\mathbf{a}^{(m)} \mid \mathbf{h}^{(n)})} \right]. \quad (12)$$

As with REINFORCE, the rewards are shifted by an offset $b(\mathbf{h})$. In this gradient, the model log-probability of a sample action sequence $\mathbf{a}^{(k)}$ is reinforced if the corresponding reward is large, or the corresponding importance weights are large, meaning that the action sequence is under-appreciated. The normalized importance weights are computed using a softmax operator $\text{softmax}(r/\tau - \log \pi_\theta)$.

5 RELATED WORK

Before presenting the experimental results, we briefly review some pieces of previous work that closely relate to the UREX approach.

Reward-Weighted Regression. Both RAML and UREX objectives bear some similarity to a method in continuous control known as Reward-Weighted Regression (RWR) (Peters & Schaal, 2007; Wierstra et al., 2008). Using our notation, the RWR objective is expressed as,

$$\mathcal{O}_{\text{RWR}}(\boldsymbol{\theta}; \tau \mid \mathbf{h}) = \log \sum_{\mathbf{a} \in \mathcal{A}} \pi_\tau^*(\mathbf{a} \mid \mathbf{h}) \pi_\theta(\mathbf{a} \mid \mathbf{h}) \quad (13)$$

$$\geq \sum_{\mathbf{a} \in \mathcal{A}} q(\mathbf{a} \mid \mathbf{h}) \log \frac{\pi_\tau^*(\mathbf{a} \mid \mathbf{h}) \pi_\theta(\mathbf{a} \mid \mathbf{h})}{q(\mathbf{a} \mid \mathbf{h})}. \quad (14)$$

To optimize \mathcal{O}_{RWR} , Peters & Schaal (2007) propose a technique inspired by the EM algorithm to maximize a variational lower bound in (14) based on a variational distribution $q(\mathbf{a} \mid \mathbf{h})$. The RWR objective can be interpreted as a log of the correlation between π_τ^* and π_θ . By contrast, the RAML and UREX objectives are both based on a KL divergence between π_τ^* and π_θ .

To optimize the RWR objective, one formulates the gradient as,

$$\frac{d}{d\theta} \mathcal{O}_{\text{RWR}}(\theta; \tau \mid \mathbf{h}) = \sum_{\mathbf{a} \in \mathcal{A}} \frac{\pi_\tau^*(\mathbf{a} \mid \mathbf{h}) \pi_\theta(\mathbf{a} \mid \mathbf{h})}{C} \frac{d}{d\theta} \log \pi_\theta(\mathbf{a} \mid \mathbf{h}), \quad (15)$$

where C denotes the normalization factor, *i.e.*, $C = \sum_{\mathbf{a} \in \mathcal{A}} \pi_\tau^*(\mathbf{a} \mid \mathbf{h}) \pi_\theta(\mathbf{a} \mid \mathbf{h})$. The expectation with respect to $\pi_\tau^*(\mathbf{a} \mid \mathbf{h}) \pi_\theta(\mathbf{a} \mid \mathbf{h}) / C$ on the RHS can be approximated by self-normalized importance sampling,¹ where the proposal distribution is π_θ . Accordingly, one draws K Monte Carlo samples $\{\mathbf{a}^{(k)}\}_{k=1}^K$ *i.i.d.* from $\pi_\theta(\mathbf{a} \mid \mathbf{h})$ and formulates the gradient as,

$$\frac{d}{d\theta} \mathcal{O}_{\text{RWR}}(\theta; \tau \mid \mathbf{h}) \approx \frac{1}{K} \sum_{k=1}^K \frac{u(\mathbf{a}^{(k)} \mid \mathbf{h})}{\sum_{m=1}^K u(\mathbf{a}^{(m)} \mid \mathbf{h})} \frac{d}{d\theta} \log \pi_\theta(\mathbf{a}^{(k)} \mid \mathbf{h}), \quad (16)$$

where $u(\mathbf{a}^{(k)} \mid \mathbf{h}) = \exp\{\frac{1}{\tau} r(\mathbf{a}^{(k)} \mid \mathbf{h})\}$. There is some similarity between (16) and (9) in that they both use self-normalized importance sampling, but note the critical difference that (16) and (9) estimate the gradients of two different objectives, and hence the importance weights in (16) do not correct for the sampling distribution $\pi_\theta(\mathbf{a} \mid \mathbf{h})$ as opposed to (9).

Beyond important technical differences, the optimal policy of \mathcal{O}_{RWR} is a one hot distribution with all probability mass concentrated on an action sequence with maximal reward, whereas the optimal policies for RAML and UREX are everywhere nonzero, with the probability of different action sequences being assigned proportionally to their exponentiated reward (with UREX introducing an additional re-scaling; see Appendix A). Further, the notion of under-appreciated reward exploration evident in $\mathcal{O}_{\text{UREX}}$, which is key to UREX’s performance, is missing in the RWR formulation.

Exploration. The RL literature contains many different attempts at incorporating exploration that may be compared with our method. The most common exploration strategy considered in value-based RL is ϵ -greedy Q-learning, where at each step the agent either takes the best action according to its current value approximation or with probability ϵ takes an action sampled uniformly at random. Like entropy regularization, such an approach applies undirected exploration, but it has achieved recent success in game playing environments (Mnih et al., 2013; Van Hasselt et al., 2016; Mnih et al., 2016).

Prominent approaches to improving exploration beyond ϵ -greedy in value-based or model-based RL have focused on reducing uncertainty by prioritizing exploration toward states and actions where the agent knows the least. This basic intuition underlies work on counter and recency methods (Thrun, 1992), exploration methods based on uncertainty estimates of values (Kaelbling, 1993; Tokic, 2010), methods that prioritize learning environment dynamics (Kearns & Singh, 2002; Stadie et al., 2015), and methods that provide an intrinsic motivation or curiosity bonus for exploring unknown states (Schmidhuber, 2006; Bellemare et al., 2016).

In contrast to value-based methods, exploration for policy-based RL methods is often a by-product of the optimization algorithm itself. Since algorithms like REINFORCE and Thompson sampling choose actions according to a stochastic policy, sub-optimal actions are chosen with some non-zero probability. The Q-learning algorithm may also be modified to sample an action from the softmax of the Q values rather than the argmax (Sutton & Barto, 1998).

Asynchronous training has also been reported to have an exploration effect on both value- and policy-based methods. Mnih et al. (2016) report that asynchronous training can stabilize training by reducing the bias experienced by a single trainer. By using multiple separate trainers, an agent is less likely to become trapped at a policy found to be locally optimal only due to local conditions. In the same spirit, Osband et al. (2016) use multiple Q value approximators and sample only one to act for each episode as a way to implicitly incorporate exploration.

By relating the concepts of value and policy in RL, the exploration strategy we propose tries to bridge the discrepancy between the two. In particular, UREX can be viewed as a hybrid combination of value-based and policy-based exploration strategies that attempts to capture the benefits of each.

¹Bornschein & Bengio (2014) apply the same trick to optimize the log-likelihood of latent variable models.

Per-step Reward. Finally, while we restrict ourselves to episodic settings where a reward is associated with an entire episode of states and actions, much work has been done to take advantage of environments that provide per-step rewards. These include policy-based methods such as actor-critic (Mnih et al., 2016; Schulman et al., 2016) and value-based approaches based on Q-learning (Van Hasselt et al., 2016; Schaul et al., 2016). Some of these value-based methods have proposed a softening of Q-values which can be interpreted as adding a form of maximum-entropy regularizer (Asadi & Littman, 2016; Azar et al., 2012; Fox et al., 2016; Ziebart, 2010). The episodic total-reward setting that we consider is naturally harder since the credit assignment to individual actions within an episode is unclear.

6 SIX ALGORITHMIC TASKS

We assess the effectiveness of the proposed approach on five algorithmic tasks from the OpenAI Gym (Brockman et al., 2016), as well as a new *binary search* problem. Each task is summarized below, with further details available on the Gym website² or in the corresponding open-source code.³ In each case, the environment has a hidden tape and a hidden sequence. The agent observes the sequence via a pointer to a single character, which can be moved by a set of *pointer control actions*. Thus an action \mathbf{a}_t is represented as a tuple (m, w, o) where m denotes how to move, w is a boolean denoting whether to write, and o is the output symbol to write.

1. **Copy:** The agent should emit a copy of the sequence. The pointer actions are move left and right.
2. **DuplicatedInput:** In the hidden tape, each character is repeated twice. The agent must deduplicate the sequence and emit every other character. The pointer actions are move left and right.
3. **RepeatCopy:** The agent should emit the hidden sequence once, then emit the sequence in the reverse order, then emit the original sequence again. The pointer actions are move left and right.
4. **Reverse:** The agent should emit the hidden sequence in the reverse order. As before, the pointer actions are move left and right.
5. **ReversedAddition:** The hidden tape is a $2 \times n$ grid of digits representing two numbers in base 3 in little-endian order. The agent must emit the sum of the two numbers, in little-endian order. The allowed pointer actions are move left, right, up, or down.

The OpenAI Gym provides an additional harder task called ReversedAddition3, which involves adding three numbers. We omit this task, since none of the methods make much progress on it.

For these tasks, the input sequences encountered during training range from a length of 2 to 33 characters. A reward of 1 is given for each correct emission. On an incorrect emission, a small penalty of -0.5 is incurred and the episode is terminated. The agent is also terminated and penalized with a reward of -1 if the episode exceeds a certain number of steps. For the experiments using UREX and MENT, we associate an episodic sequence of actions with the *total* reward, defined as the sum of the per-step rewards. The experiments using Q-learning, on the other hand, used the per-step rewards. Each of the Gym tasks has a *success threshold*, which determines the required average reward over 100 episodes for the agent to be considered successful.

We also conduct experiments on an additional algorithmic task described below:

6. **BinarySearch:** Given an integer n , the environment has a hidden array of n distinct numbers stored in ascending order. The environment also has a query number x unknown to the agent that is contained somewhere in the array. The goal of the agent is to find the query number in the array in a small number of actions. The environment has three integer registers initialized at $(n, 0, 0)$. At each step, the agent can interact with the environment via the four following actions:
 - $\text{INC}(i)$: increment the value of the register i for $i \in \{1, 2, 3\}$.
 - $\text{DIV}(i)$: divide the value of the register i by 2 for $i \in \{1, 2, 3\}$.
 - $\text{AVG}(i)$: replace the value of the register i with the average of the two other registers.
 - $\text{CMP}(i)$: compare the value of the register i with x and receive a signal indicating which value is greater. The agent succeeds when it calls CMP on an array cell holding the value x .

²gym.openai.com

³github.com/openai/gym

Table 1: Each cell shows the percentage of 60 trials with different hyper-parameters (η , c) and random restarts that successfully solve an algorithmic task. UREX is more robust to hyper-parameter changes than MENT. We evaluate MENT with a few temperatures and UREX with $\tau = 0.1$.

	REINFORCE / MENT				UREX
	$\tau = 0.0$	$\tau = 0.005$	$\tau = 0.01$	$\tau = 0.1$	$\tau = 0.1$
Copy	85.0	88.3	90.0	3.3	75.0
DuplicatedInput	68.3	73.3	73.3	0.0	100.0
RepeatCopy	0.0	0.0	11.6	0.0	18.3
Reverse	0.0	0.0	3.3	10.0	16.6
ReversedAddition	0.0	0.0	1.6	0.0	30.0
BinarySearch	0.0	0.0	1.6	0.0	20.0

The agent is terminated when the number of steps exceeds a maximum threshold of $2n+1$ steps and receives a reward of 0. If the agent finds x at step t , it receives a reward of $10(1-t/(2n+1))$.

We set the maximum number of steps to $2n+1$ to allow the agent to perform a full linear search. A policy performing full linear search achieves an average reward of 5, because x is chosen uniformly at random from the elements of the array. A policy employing binary search can find the number x in at most $2\log_2 n + 1$ steps. If n is selected uniformly at random from the range $32 \leq n \leq 512$, binary search yields an optimal average reward above 9.55. We set the *success threshold* for this task to an average reward of 9.

7 EXPERIMENTS

We compare our policy gradient method using under-appreciated reward exploration (UREX) against two main RL baselines: (1) REINFORCE with entropy regularization termed MENT (Williams & Peng, 1991), where the value of τ determines the degree of regularization. When $\tau = 0$, standard REINFORCE is obtained. (2) one-step double Q-learning based on bootstrapping one step future rewards.

7.1 ROBUSTNESS TO HYPER-PARAMETERS

Hyper-parameter tuning is often tedious for RL algorithms. We found that the proposed UREX method significantly improves robustness to changes in hyper-parameters when compared to MENT. For our experiments, we perform a careful grid search over a set of hyper-parameters for both MENT and UREX. For any hyper-parameter setting, we run the MENT and UREX methods 5 times with different random restarts. We explore the following main hyper-parameters:

- The *learning rate* denoted η chosen from a set of 3 possible values $\eta \in \{0.1, 0.01, 0.001\}$.
- The maximum L2 norm of the gradients, beyond which the gradients are *clipped*. This parameter, denoted c , matters for training RNNs. The value of c is selected from $c \in \{1, 10, 40, 100\}$.
- The temperature parameter τ that controls the degree of exploration for both MENT and UREX. For MENT, we use $\tau \in \{0, 0.005, 0.01, 0.1\}$. For UREX, we only consider $\tau = 0.1$, which consistently performs well across the tasks.

In all of the experiments, both MENT and UREX are treated exactly the same. In fact, the change of implementation is just a few lines of code. Given a value of τ , for each task, we run 60 training jobs comprising 3 learning rates, 4 clipping values, and 5 random restarts. We run each algorithm for a maximum number of steps determined based on the difficulty of the task. The training jobs for Copy, DuplicatedInput, RepeatCopy, Reverse, ReversedAddition, and BinarySearch are run for $2K$, 500, $50K$, $5K$, $50K$, and $2K$ stochastic gradient steps, respectively. We find that running a trainer job longer does not result in a better performance. Our policy network comprises a single LSTM layer with 128 nodes. We use the Adam optimizer (Kingma & Ba, 2015) for the experiments.

Table 1 shows the percentage of 60 trials on different hyper-parameters (η , c) and random restarts which successfully solve each of the algorithmic tasks. It is clear that UREX is more robust than

MENT to changes in hyper-parameters, even though we only report the results of UREX for a single temperature. See Appendix B for more detailed tables on hyper-parameter robustness.

7.2 RESULTS

Table 2 presents the number of successful attempts (out of 5 random restarts) and the expected reward values (averaged over 5 trials) for each RL algorithm given the best hyper-parameters. One-step Q-learning results are also included in the table. We also present the training curves for MENT and UREX in Figure 2. It is clear that UREX outperforms the baselines on these tasks. On the more difficult tasks, such as Reverse and ReverseAddition, UREX is able to consistently find an appropriate algorithm, but MENT and Q-learning fall behind. Importantly, for the BinarySearch task, which exhibits many local maxima and necessitates smart exploration, UREX is the only method that can solve it consistently. The Q-learning baseline solves some of the simple tasks, but it makes little headway on the harder tasks. We believe that entropy regularization for policy gradient and ϵ -greedy for Q-learning are relatively weak exploration strategies in long episodic tasks with delayed rewards. On such tasks, one random exploratory step in the wrong direction can take the agent off the optimal policy, hampering its ability to learn. In contrast, UREX provides a form of adaptive and smart exploration. In fact, we observe that the variance of the importance weights decreases as the agent approaches the optimal policy, effectively reducing exploration when it is no longer necessary; see Appendix E.

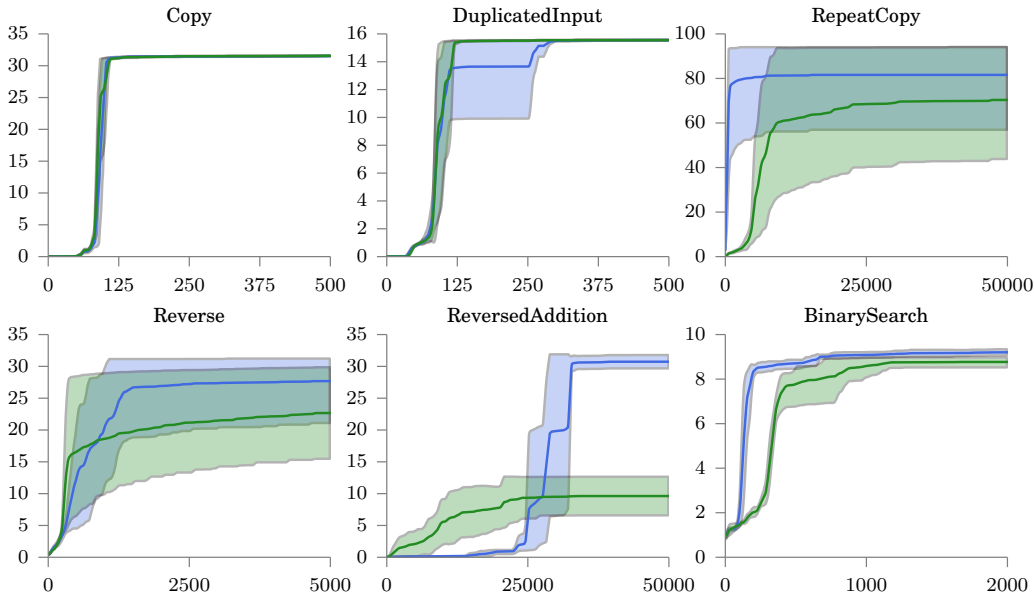


Figure 2: Average reward during training for MENT (green) and UREX (blue). We find the best hyper-parameters for each method, and run each algorithm 5 times with random restarts. The curves present the average reward as well as the single standard deviation region clipped at the min and max.

7.3 GENERALIZATION TO LONGER SEQUENCES

To confirm whether our method is able to find the correct algorithm for multi-digit addition, we investigate its generalization to longer input sequences than provided during training. We evaluate the trained models on inputs up to a length of 2000 digits, even though training sequences were at most 33 characters. For each length, we test the model on 100 randomly generated inputs, stopping when the accuracy falls below 100%. Out of the 60 models trained on addition with UREX, we find that 5 models generalize to numbers up to 2000 digits without any observed mistakes. On the best UREX hyper-parameters, 2 out of the 5 random restarts are able to generalize successfully. For more detailed results on the generalization performance on 3 different tasks including Copy,

Table 2: Results on several algorithmic tasks comparing Q-learning and policy gradient based on MENT and UREX. We find the best hyper-parameters for each method, and run each algorithm 5 times with random restarts. Number of successful attempts (out of 5) that achieve a reward threshold is reported. Expected reward computed over the last few iterations of training is also reported.

	Num. of successful attempts out of 5			Expected reward		
	Q-learning	MENT	UREX	Q-learning	MENT	UREX
Copy	5	5	5	31.2	31.2	31.2
DuplicatedInput	5	5	5	15.4	15.4	15.4
RepeatCopy	1	3	4	39.3	69.2	81.1
Reverse	0	2	4	4.4	21.9	27.2
ReversedAddition	0	1	5	1.1	8.7	30.2
BinarySearch	0	1	4	5.2	8.6	9.1

DuplicatedInput, and ReversedAddition, see Appendix C. During these evaluations, we take the action with largest probability from $\pi_\theta(\mathbf{a} \mid \mathbf{h})$ at each time step rather than sampling randomly.

We also looked into the generalization of the models trained on the BinarySearch task. We found that none of the agents perform proper binary search. Rather, those that solved the task perform a hybrid of binary and linear search: first actions follow a binary search pattern, but then the agent switches to a linear search procedure once it narrows down the search space; see Appendix D for some execution traces for BinarySearch and ReversedAddition. Thus, on longer input sequences, the agent’s running time complexity approaches linear rather than logarithmic. We hope that future work will make more progress on this task. This task is especially interesting because the reward signal should incorporate both correctness and efficiency of the algorithm.

7.4 IMPLEMENTATION DETAILS

In all of the experiments, we make use of curriculum learning. The environment begins by only providing small inputs and moves on to longer sequences once the agent achieves close to maximal reward over a number of steps. For policy gradient methods including MENT and UREX, we only provide the agent with a reward at the end of the episode, and there is no notion of intermediate reward. For the value-based baseline, we implement one-step Q-learning as described in Mnih et al. (2016)-Alg. 1, employing double Q-learning with ϵ -greedy exploration. We use the same RNN in our policy-based approaches to estimate the Q values. A grid search over exploration rate, exploration rate decay, learning rate, and sync frequency (between online and target network) is conducted to find the best hyper-parameters. Unlike our other methods, the Q-learning baseline uses intermediate rewards, as given by the OpenAI Gym on a per-step basis. Hence, the Q-learning baseline has a slight advantage over the policy gradient methods.

In all of the tasks except Copy, our stochastic optimizer uses mini-batches comprising 400 policy samples from the model. These 400 samples correspond to 40 different random sequences drawn from the environment, and 10 random policy trajectories per sequence. In other words, we set $K = 10$ and $N = 40$ as defined in (3) and (12). For MENT, we use the 10 samples to subtract the mean of the coefficient of $\frac{d}{d\theta} \log \pi_\theta(\mathbf{a} \mid \mathbf{h})$ which includes the contribution of the reward and entropy regularization. For UREX, we use the 10 trajectories to subtract the mean reward and normalize the importance sampling weights. We do not subtract the mean of the normalized importance weights. For the Copy task, we use mini-batches with 200 samples using $K = 10$ and $N = 20$. Experiments are conducted using Tensorflow (Abadi et al., 2016).

8 CONCLUSION

We present a variant of policy gradient, called UREX, which promotes the exploration of action sequences that yield rewards larger than what the model expects. This exploration strategy is the result of importance sampling from the optimal policy. Our experimental results demonstrate that UREX significantly outperforms other value and policy based methods, while being more robust

to changes of hyper-parameters. By using UREX, we can solve algorithmic tasks like multi-digit addition from only episodic reward, which other methods cannot reliably solve even given the best hyper-parameters. We introduce a new algorithmic task based on binary search to advocate more research in this area, especially when the computational complexity of the solution is also of interest. Solving these tasks is not only important for developing more human-like intelligence in learning algorithms, but also important for generic reinforcement learning, where smart and efficient exploration is the key to successful methods.

9 ACKNOWLEDGMENT

We thank Sergey Levine, Irwan Bello, Corey Lynch, George Tucker, Kelvin Xu, Volodymyr Mnih, and the Google Brain team for insightful comments and discussions.

REFERENCES

- Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, et al. Tensorflow: A system for large-scale machine learning. *arXiv:1605.08695*, 2016.
- Dana Angulin. Learning regular sets from queries and counterexamples. *Information and Computation*, 1987.
- Kavosh Asadi and Michael L Littman. A new softmax operator for reinforcement learning. *arXiv preprint arXiv:1612.05628*, 2016.
- Mohammad Gheshlaghi Azar, Vicenç Gómez, and Hilbert J Kappen. Dynamic policy programming. *Journal of Machine Learning Research*, 13(Nov):3207–3245, 2012.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *ICLR*, 2015.
- Marc G. Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. *NIPS*, 2016.
- Jörg Bornschein and Yoshua Bengio. Reweighted wake-sleep. *arXiv:1406.2751*, 2014.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. OpenAI Gym. *arXiv:1606.01540*, 2016.
- Roy Fox, Ari Pakman, and Naftali Tishby. G-learning: Taming the noise in reinforcement learning via soft updates. *Uncertainty in Artificial Intelligence*, 2016. URL <http://arxiv.org/abs/1512.08562>.
- Gene Golub. Some modified matrix eigenvalue problems. *SIAM Review*, 1987.
- Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwinska, Sergio G. Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, Adria P. Badia, Karl M. Hermann, Yori Zwols, Georg Ostrovski, Adam Cain, Helen King, Christopher Summerfield, Phil Blunsom, Koray Kavukcuoglu, and Demis Hassabis. Hybrid computing using a neural network with dynamic external memory. *Nature*, 2016.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 1997.
- Leslie Pack Kaelbling. *Learning in embedded systems*. MIT press, 1993.
- Lukasz Kaiser and Ilya Sutskever. Neural GPUs learn algorithms. *ICLR*, 2016.
- Michael Kearns and Satinder Singh. Near-optimal reinforcement learning in polynomial time. *Machine Learning*, 2002.
- Charles Kemp, Noah Goodman, and Joshua Tenenbaum. Learning and using relational theories. *NIPS*, 2007.

- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- N. Lavrac and S. Dzeroski. *Inductive Logic Programming: Theory and Methods*. Ellis Horwood, 1994.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *arXiv:1312.5602*, 2013.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, et al. Human-level control through deep reinforcement learning. *Nature*, 2015.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy P Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. *ICML*, 2016.
- Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT Press, 2012.
- Arvind Neelakantan, Quoc V. Le, and Ilya Sutskever. Neural programmer: Inducing latent programs with gradient descent. *ICLR*, 2016.
- Mohammad Norouzi, Samy Bengio, Zhifeng Chen, Navdeep Jaitly, Mike Schuster, Yonghui Wu, and Dale Schuurmans. Reward augmented maximum likelihood for neural structured prediction. *NIPS*, 2016.
- Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped DQN. *NIPS*, 2016.
- Art B. Owen. *Monte Carlo theory, methods and examples*. 2013.
- Jan Peters and Stefan Schaal. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pp. 745–750. ACM, 2007.
- Scott E. Reed and Nando de Freitas. Neural programmer-interpreters. *ICLR*, 2016.
- Stuart Jonathan Russell, Peter Norvig, John F Canny, Jitendra M Malik, and Douglas D Edwards. *Artificial intelligence: a modern approach*, volume 2. Prentice hall Upper Saddle River, 2003.
- Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *ICLR*, 2016.
- Jürgen Schmidhuber. Optimal artificial curiosity, creativity, music, and the fine arts. *Connection Science*, 2006.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *ICLR*, 2016.
- David Silver, Aja Huang, et al. Mastering the game of Go with deep neural networks and tree search. *Nature*, 2016.
- Bradly C. Stadie, Sergey Levine, and Pieter Abbeel. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv:1507.00814*, 2015.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. *NIPS*, 2014.
- Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT Press, 1998.
- Sebastian B Thrun. Efficient exploration in reinforcement learning. Technical report, 1992.
- Michel Tokic. Adaptive ϵ -greedy exploration in reinforcement learning based on value differences. *AAAI*, 2010.
- Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. *AAAI*, 2016.

- Oriol Vinyals, Meire Fortunato, and Navdeep Jaitly. Pointer networks. *NIPS*, 2015.
- Daan Wierstra, Tom Schaul, Jan Peters, and Juergen Schmidhuber. Episodic reinforcement learning by logistic reward-weighted regression. In *International Conference on Artificial Neural Networks*, pp. 407–416. Springer, 2008.
- Ronald J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 1992.
- Ronald J Williams and Jing Peng. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 1991.
- Wojciech Zaremba and Ilya Sutskever. Learning to execute. *arXiv:1410.4615*, 2014.
- Wojciech Zaremba and Ilya Sutskever. Reinforcement learning neural turing machines. *arXiv:1505.00521*, 2015.
- Brian D Ziebart. Modeling purposeful adaptive behavior with the principle of maximum causal entropy. 2010.

A OPTIMAL POLICY FOR THE UREX OBJECTIVE

To derive the form of the optimal policy for the UREX objective (11), note that for each \mathbf{h} one would like to maximize

$$\sum_{\mathbf{a} \in \mathcal{A}} \left[\pi_{\theta}(\mathbf{a}) r(\mathbf{a}) + \tau \pi_{\tau}^*(\mathbf{a}) \log \pi_{\theta}(\mathbf{a}) \right], \quad (17)$$

subject to the constraint $\sum_{\mathbf{a} \in \mathcal{A}} \pi_{\theta}(\mathbf{a}) = 1$. To enforce the constraint, we introduce a Lagrange multiplier α and aim to maximize

$$\sum_{\mathbf{a} \in \mathcal{A}} \left[\pi_{\theta}(\mathbf{a}) r(\mathbf{a}) + \tau \pi_{\tau}^*(\mathbf{a}) \log \pi_{\theta}(\mathbf{a}) - \alpha \pi_{\theta}(\mathbf{a}) \right] + \alpha. \quad (18)$$

Since the gradient of the Lagrangian (18) with respect to θ is given by

$$\sum_{\mathbf{a} \in \mathcal{A}} \frac{d\pi_{\theta}(\mathbf{a})}{d\theta} \left[r(\mathbf{a}) + \tau \frac{\pi_{\tau}^*(\mathbf{a})}{\pi_{\theta}(\mathbf{a})} - \alpha \right], \quad (19)$$

the optimal choice for π_{θ} is achieved by setting

$$\pi_{\theta}(\mathbf{a}) = \frac{\tau \pi_{\tau}^*(\mathbf{a})}{\alpha - r(\mathbf{a})} \quad \text{for all } \mathbf{a} \in \mathcal{A}, \quad (20)$$

forcing the gradient to be zero. The Lagrange multiplier α can then be chosen so that $\sum_{\mathbf{a} \in \mathcal{A}} \pi_{\theta}(\mathbf{a}) = 1$ while also satisfying $\alpha > \max_{\mathbf{a} \in \mathcal{A}} r(\mathbf{a})$; see *e.g.* (Golub, 1987).

B ROBUSTNESS TO HYPER-PARAMETERS

Tables 3–8 provide more details on different cells of Table 1. Each table presents the results of MENT using the best temperature τ vs. UREX with $\tau = 0.1$ on a variety of learning rates and clipping values. Each cell is the number of trials out of 5 random restarts that succeed at solving the task using a specific η and c .

Table 3: Copy – number of successful attempts out of 5.

	MENT ($\tau = 0.01$)			UREX ($\tau = 0.1$)		
	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$
$c = 1$	3	5	5	5	5	2
$c = 10$	5	4	5	5	5	3
$c = 40$	3	5	5	4	4	1
$c = 100$	4	5	5	4	5	2

Table 4: DuplicatedInput – number of successful attempts out of 5.

	MENT ($\tau = 0.01$)			UREX ($\tau = 0.1$)		
	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$
$c = 1$	3	5	3	5	5	5
$c = 10$	2	5	3	5	5	5
$c = 40$	4	5	3	5	5	5
$c = 100$	2	5	4	5	5	5

Table 5: RepeatCopy – number of successful attempts out of 5.

	MENT ($\tau = 0.01$)			UREX ($\tau = 0.1$)		
	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$
$c = 1$	0	1	0	0	2	0
$c = 10$	0	0	2	0	4	0
$c = 40$	0	0	1	0	2	0
$c = 100$	0	0	3	0	3	0

Table 6: Reverse – number of successful attempts out of 5.

	MENT ($\tau = 0.1$)			UREX ($\tau = 0.1$)		
	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$
$c = 1$	1	1	0	0	0	0
$c = 10$	0	1	0	0	4	0
$c = 40$	0	2	0	0	2	1
$c = 100$	1	0	0	0	2	1

Table 7: ReversedAddition – number of successful attempts out of 5.

	MENT ($\tau = 0.01$)			UREX ($\tau = 0.1$)		
	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$
$c = 1$	0	0	0	0	0	4
$c = 10$	0	0	0	0	3	2
$c = 40$	0	0	0	0	0	5
$c = 100$	0	0	1	0	1	3

Table 8: BinarySearch – number of successful attempts out of 5.

	MENT ($\tau = 0.01$)			UREX ($\tau = 0.1$)		
	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$	$\eta = 0.1$	$\eta = 0.01$	$\eta = 0.001$
$c = 1$	0	0	0	0	4	0
$c = 10$	0	1	0	0	3	0
$c = 40$	0	0	0	0	3	0
$c = 100$	0	0	0	0	2	0

C GENERALIZATION TO LONGER SEQUENCES

Table 9 provides a more detailed look into the generalization performance of the trained models on Copy, DuplicatedInput, and ReversedAddition. The tables show how the number of models which can solve the task correctly drops off as the length of the input increases.

Table 9: Generalization results. Each cell includes the number of runs out of 60 different hyperparameters and random initializations that achieve 100% accuracy on input of length up to the specified length. The bottom row is the maximal length (≤ 2000) up to which at least one model achieves 100% accuracy.

	Copy		DuplicatedInput		ReversedAddition	
	MENT	UREX	MENT	UREX	MENT	UREX
30	54	45	44	60	1	18
100	51	45	36	56	0	6
500	27	22	19	25	0	5
1000	3	2	12	17	0	5
2000	0	0	6	9	0	5
Max	1126	1326	2000	2000	38	2000

D EXAMPLE EXECUTION TRACES

We provide the traces of two trained agents on the ReversedAddition task (Figure 3) and the BinarySearch task (Table 10).

E VARIANCE OF IMPORTANCE WEIGHTS

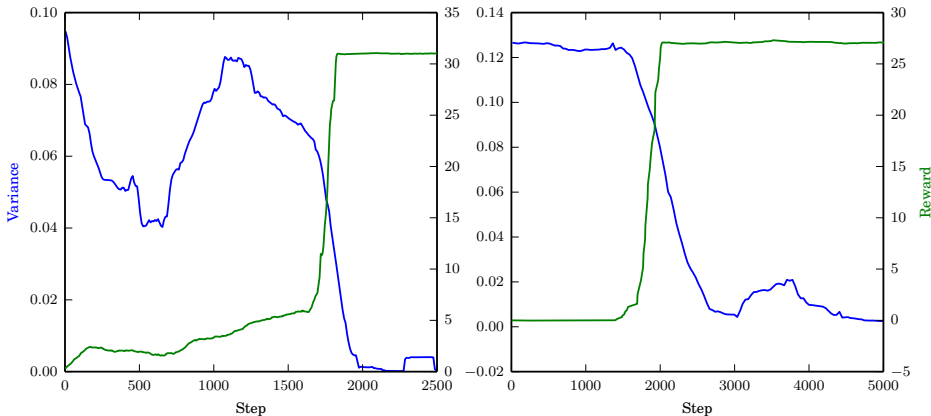


Figure 4: This plot shows the variance of the importance weights in the UREX updates as well as the average reward for two successful runs. We see that the variance starts off high and reaches near zero towards the end when the optimal policy is found. In the first plot, we see a dip and rise in the variance which corresponds to a plateau and then increase in the average reward.

F A SIMPLE BANDIT TASK

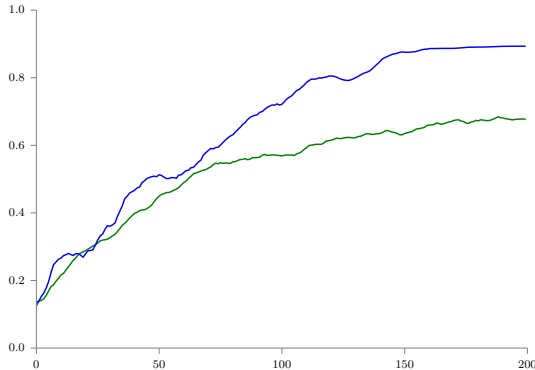


Figure 5: In this plot we present the average performance of UREX (blue) and MENT (green) over 100 repeats of a bandit-like task after choosing optimal hyperparameters for each method. In the task, the agent chooses one of 10,000 actions at each step and receives a payoff corresponding to the entry in a reward vector $\mathbf{r} = (r_1, \dots, r_{10,000})$ such that $r_i = u_i^\beta$, where $u_i \in [0, 1]$ has been sampled randomly and independently from a uniform distribution. We parameterize the policy with a weight vector $\theta \in \mathbb{R}^{30}$ such that $\pi_\theta(a) \propto \exp(\phi^{(a)} \cdot \theta)$, where the basis vectors $\phi^{(a)} \in \mathbb{R}^{30}$ for each action are sampled from a standard normal distribution. The plot shows the average rewards obtained by setting $\beta = 8$ over 100 experiments, consisting of 10 repeats (where \mathbf{r} and $\Phi = (\phi^{(1)}, \dots, \phi^{(10,000)})$ are redrawn at the start of each repeat), with 10 random restarts within each repeat (keeping \mathbf{r} and Φ fixed but reinitializing θ). Thus, this task presents a relatively simple problem with a large action space, and we again see that UREX outperforms MENT.