

Learning by Competition of Self-Interested Reinforcement Learning Agents

Stephen Chung

Department of Computer Science, University of Massachusetts Amherst, USA
minghaychung@umass.edu

Abstract

An artificial neural network can be trained by uniformly broadcasting a reward signal to units that implement a REINFORCE learning rule. Though this presents a biologically plausible alternative to backpropagation in training a network, the high variance associated with it renders it impractical to train deep networks. The high variance arises from the inefficient structural credit assignment since a single reward signal is used to evaluate the collective action of all units. To facilitate structural credit assignment, we propose replacing the reward signal to hidden units with the change in the L^2 norm of the unit's outgoing weight. As such, each hidden unit in the network is trying to maximize the norm of its outgoing weight instead of the global reward, and thus we call this learning method *Weight Maximization*. We prove that Weight Maximization is approximately following the gradient of rewards in expectation. In contrast to backpropagation, Weight Maximization can be used to train both continuous-valued and discrete-valued units. Moreover, Weight Maximization solves several major issues of backpropagation relating to biological plausibility. Our experiments show that a network trained with Weight Maximization can learn significantly faster than REINFORCE and slightly slower than backpropagation. Weight Maximization illustrates an example of cooperative behavior automatically arising from a population of self-interested agents in a competitive game without any central coordination.

Introduction

The *error backpropagation algorithm* (backprop) (Rumelhart, Hinton, and Williams 1986) efficiently computes the gradient of an objective function with respect to parameters by iterating backward from the last layer of a multi-layer artificial neural network (ANN). However, backprop is generally regarded as being biologically implausible (Crick 1989; Mazzoni, Andersen, and Jordan 1991; O'Reilly 1996; Bengio et al. 2015; Hassabis et al. 2017; Lillicrap et al. 2020). First, the learning rule given by backprop is non-local, as it relies on information other than input and output of a neuron-like unit computed in the feedforward phase. Second, backprop requires synaptic symmetry in the forward and backward paths, which has not been observed in biological systems. Third, backprop requires precise coordination

between the feedforward and feedback phase because the feedforward value has to be retained until the error signal arrives.

Alternatively, REINFORCE (Williams 1992), a special case of $A_{R-\lambda P}$ when $\lambda = 0$ (Barto and Anandan 1985), could be applied to all units as a more biologically plausible way of training a network. It is shown that the learning rule gives an unbiased estimate of the gradient of return (Williams 1992). Another interpretation of this relates to viewing each unit as a reinforcement learning (RL) agent, with each agent trying to maximize the same reward from the environment. We can thus view an ANN as a *team of agents* playing a cooperative game, a scenario where all agents receive the same reward. Such a team of agents is also known as *coagent network* (Thomas 2011). However, coagent networks can only feasibly solve simple tasks due to the high variance associated with this training method and thus the low speed of learning. The high variance stems from the lack of structural credit assignment, i.e. a single reward signal is used to evaluate the collective action of all agents.

To address the lack of structural credit assignment in a team of agents trained by REINFORCE, we consider delivering a different reward signal to each hidden agent instead of the same global reward. Here hidden agents correspond to hidden units in an ANN and refer to the agents that output to other agents in the team instead of to the environment. As such, each hidden agent is associated with an *outgoing weight*, that is, the vector of the weight by which the agent's actions influence other agents in the team. We propose to replace the global reward signal to each hidden agent with the change in the L^2 norm of its outgoing weight, such that each hidden agent in the team is trying to maximize the norm of its outgoing weight. We call this new learning method *Weight Maximization*. This is based on the intuition that the norm of an agent's outgoing weight roughly reflects the contribution of the agent in the team. This change of reward signals turns the original cooperative game into a competitive game since agents no longer receive the same reward.

We prove that Weight Maximization is approximately following the gradient of return in expectation, showing that every hidden agent maximizing the norm of its outgoing weight also approximately maximizes the team's rewards. This illustrates an example of cooperative behavior automatically arising from a population of self-interested agents in

a competitive game and offers an alternative perspective of training an ANN - each unit maximizing the norm of its outgoing weight instead of a network maximizing its performance. This alternative perspective localizes the optimization problem for each unit, yielding a wide range of RL solutions in training a network. Our experiments show that a network trained with Weight Maximization can learn much faster than REINFORCE, such that its speed is slightly lower than backprop.

One may question whether the change of synaptic strength, analogous to the change of weights in ANNs, can be used to guide plasticity in biological systems, since the change of synaptic strength has a much slower timescale compared to the activation of neurons. To address this issue, we generalize Weight Maximization to use eligibility traces, such that the network can still learn when the outgoing weights change slowly. Weight Maximization with eligibility traces also solves the three aforementioned problems of backprop regarding biological plausibility. Nonetheless, the biological plausibility of Weight Maximization remains to be investigated. It is not yet clear if there exists a molecular mechanism that uses the change of synaptic strength in axons to guide the change of synaptic strength in dendrites.

In summary, our paper has the following main contributions:

- We propose a novel algorithm called Weight Maximization that allows efficient structural credit assignment and significantly lowers the variance associated with REINFORCE when training a team of agents;
- We prove that Weight Maximization is approximately following the gradient of return in expectation, establishing the approximate equivalence of hidden units maximizing the norm of outgoing weights and external rewards, thus providing theoretical justification for algorithms (Uhr and Vossler 1961; Klopf and Gose 1969; Anderson 1986; Selfridge 1988) based on the norm of outgoing weights;
- Weight Maximization can be used to train continuous-valued and discrete-valued units, offering an advantage to backprop;
- Weight Maximization represents a feasible alternative to backprop given their comparable learning speed;
- We generalize Weight Maximization to use eligibility traces, which solves several major issues of backprop relating to biological plausibility.

The paper and the appendix is available at <https://arxiv.org/abs/2010.09770>.

Notation

We consider a Markov Decision Process (MDP) defined by a tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma, d_0)$, where \mathcal{S} is a finite set of states of an agent's environment (although this work can be extended to the infinite state case), \mathcal{A} is a finite set of actions, $P : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a transition function giving the dynamics of the environment, $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is a reward function, $\gamma \in [0, 1]$ is a discount factor, and $d_0 : \mathcal{S} \rightarrow [0, 1]$ is an initial state distribution. Denoting the

state, action, and reward signal at time t by S_t , A_t , and R_t respectively, $P(s, a, s') = \Pr(S_{t+1} = s' | S_t = s, A_t = a)$, $R(s, a) = \mathbb{E}[R_t | S_t = s, A_t = a]$, and $d_0(s) = \Pr(S_0 = s)$, where P and d_0 are probability mass functions. An episode is a sequence of states, actions, and rewards, starting from $t = 0$ and continuing until reaching the terminal state. For any learning methods, we can measure its performance as it improves over multiple episodes, which makes up a run.

Letting $G_t = \sum_{k=t}^{\infty} \gamma^{k-t} R_k$ denote the infinite-horizon discounted return accrued after acting at time t , we are interested in finding, or approximating, a policy $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ such that for any time $t > 0$, selecting actions according to $\pi(s, a) = \Pr(A_t = a | S_t = s)$ maximizes the expected return $\mathbb{E}[G_t]$. The value function for policy π is V^π where for all $s \in \mathcal{S}$, $V^\pi(s) = \mathbb{E}[G_t | S_t = s, \pi]$, which can be shown to be independent of t .

In this paper, we restrict attention to policies computed by a network of L stochastic units. The following definitions hold for any time $t > 0$. Let $H_t^{(l)}$ denote the activation value of the unit $l \in \{1, 2, \dots, L\}$ at time t , which is a one-dimensional discrete or continuous random variable. We also let $H_t^{(0)} = S_t$ and $A_t = H_t^{(L)}$. We call unit l , where $1 \leq l \leq L - 1$, a hidden unit and unit L the output unit. For any $1 \leq l \leq L$, the distribution of $H_t^{(l)}$ conditional on $H_t^{(0:l-1)}$ is given by $\pi_l : (\mathcal{S} \times \mathbb{R}^{l-1}) \times \mathbb{R} \rightarrow [0, 1]$, such that $\Pr(H_t^{(l)} = h^{(l)} | H_t^{(0:l-1)} = h^{(0:l-1)}; \mathbf{w}^{(l)}) = \pi_l(h^{(0:l-1)}, h^{(l)}; \mathbf{w}^{(l)})$, where $\mathbf{w}^{(l)}$ is the parameter of unit l . To sample an action A_t from the network given S_t , we iteratively sample $H_t^{(l)} \sim \pi_l(H_t^{(0:l-1)}, \cdot; \mathbf{w}^{(l)})$ from $l = 1$ to L . In other words, the network is a feedforward network with units ordered and connections restricted to higher-number units. This formulation is a generalization of multi-layer networks of stochastic units since multi-layer networks are the special case of units being arranged in layers and connections between non-adjacent layers being set to zero. We mostly focus on multi-layer networks in this paper.

We assume that for all $1 \leq l \leq L$, $\pi_l(h^{(0:l-1)}, h^{(l)}; \mathbf{w}^{(l)})$ can be expressed as a differentiable function $f_l(\mathbf{w}^{(l)} \cdot h^{(0:l-1)}, h^{(l)})$, where \cdot denotes the dot product, and call the vector $\mathbf{w}^{(l)}$ the *incoming weight* of unit l . We denote the weight connecting from unit k to l (where $k < l$) as $\mathbf{w}_{(k)}^{(l)}$, and call the vector $\mathbf{v}^{(k)} = [\mathbf{w}_{(k)}^{(k+1)}, \mathbf{w}_{(k)}^{(k+2)}, \dots, \mathbf{w}_{(k)}^{(L)}]^T$ the *outgoing weight* of unit k .

Though we restrict our attention to a network with a single output unit, the algorithm we present can be generalized to a network with multiple output units easily by only replacing the rewards to all hidden units by the change in the norm of their outgoing weights as in (4).

The case we consider here is one in which all the units of the network implement an RL algorithm and share the same reward signal. These networks can therefore be considered to be *teams of agents*, which have also been called *coagent networks* (Thomas 2011); agents here refer to *RL agents* (Sutton and Barto 2018).

We denote $\|\mathbf{x}\|_p^p$ as the p -norm of vector \mathbf{x} to the power of p , and $x^{(m:n)}$ as $\{x^m, x^{m+1}, \dots, x^n\}$.

Algorithm

The gradient of return at time t with respect to $\mathbf{w}^{(l)}$, where $1 \leq l \leq L$, can be estimated by REINFORCE, also known as the likelihood ratio estimator:

$$\nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G_t] = \sum_{k=t}^{\infty} \gamma^{(k-t)} \mathbb{E}[G_k \nabla_{\mathbf{w}^{(l)}} \log \Pr(A_k | S_k)]. \quad (1)$$

We can show that the terms in the summation of (1) also equals the expectation of the update given by REINFORCE applied to each hidden unit with the same reinforcement signal G_t :

Theorem 1. *Let the policy be a network of stochastic units as defined above. For all $t > 0$ and $1 \leq l \leq L$,*

$$\begin{aligned} & \mathbb{E}[G_t \nabla_{\mathbf{w}^{(l)}} \log \Pr(A_t | S_t)] \\ &= \mathbb{E}[G_t \nabla_{\mathbf{w}^{(l)}} \log \pi_l(H_t^{(0:l-1)}, H_t^{(l)}; \mathbf{w}^{(l)})]. \end{aligned} \quad (2)$$

See Williams (1992) for the proof. This shows that we can apply REINFORCE to each unit of the network, and the learning rule still gives an unbiased estimate of the gradient of the return. Therefore, denoting α as the step size, we can update parameters by the following stochastic gradient ascent rule:

$$\mathbf{w}^{(l)} \leftarrow \mathbf{w}^{(l)} + \alpha G_t \nabla_{\mathbf{w}^{(l)}} \log \pi_l(H_t^{(0:l-1)}, H_t^{(l)}; \mathbf{w}^{(l)}). \quad (3)$$

However, this learning rule suffers from high variance since a single reinforcement signal (G_t) is used to evaluate the collective action of all units. In other words, the signal G_t has a weak correlation with the gradient $\nabla_{\mathbf{w}^{(l)}} \log \pi_l(H_t^{(0:l-1)}, H_t^{(l)}; \mathbf{w}^{(l)})$, making the multiplication of these two terms have a high variance. To reduce this variance, we propose to replace this signal to each unit l by:

$$G_t^{(l)} = \begin{cases} \mathbf{v}^{(l)} \cdot \Delta \mathbf{v}_t^{(l)} & \text{for } l \in \{1, 2, \dots, L-1\}, \\ G_t & \text{for } l = L. \end{cases} \quad (4)$$

where \cdot denotes the dot product, and $\Delta \mathbf{v}_t^{(l)}$ is the change of the outgoing weight, $\mathbf{v}^{(l)}$, resulting from the update of the outgoing weight at time t . For the output unit, we let $G_t^{(L)} = G_t$; that is, the output unit is still maximizing the return from the environment.

With the new reinforcement signal $G_t^{(l)}$, each hidden unit is approximately maximizing the L^2 norm of its outgoing weight. To see this, consider the change in the L^2 norm of the outgoing weight for hidden unit l at time t :

$$\|\mathbf{v}^{(l)} + \Delta \mathbf{v}_t^{(l)}\|_2^2 - \|\mathbf{v}^{(l)}\|_2^2 \quad (5)$$

$$= 2\mathbf{v}^{(l)} \cdot \Delta \mathbf{v}_t^{(l)} + \|\Delta \mathbf{v}_t^{(l)}\|_2^2. \quad (6)$$

We observe that $G_t^{(l)}$ is proportional to (6) except for the term $\|\Delta \mathbf{v}_t^{(l)}\|_2^2$. We choose to ignore this term in the reinforcement signal since this term is $\mathcal{O}(\alpha^2)$ while $2\mathbf{v}^{(l)} \cdot \Delta \mathbf{v}_t^{(l)} = \mathcal{O}(\alpha)$. If the step size α is very small as in typical experiments, then $\|\Delta \mathbf{v}_t^{(l)}\|_2^2$ is also negligible and does not affect experimental results. However, by adjusting α , this

term can be made arbitrarily small or large, and so we choose to remove it completely instead. This removal is necessary to arrive Theorem 2.

The motivation of using the change in the norm of a unit's outgoing weight as a reinforcement signal in (4) is based on the idea that the norm of a unit's outgoing weight roughly reflects the contribution of the unit in the network. For example, if the hidden unit's output is useful in guiding action, then the output unit will learn a large weight associated with it. Conversely, if the hidden agent is outputting random noise, then the output unit will learn a zero weight associated with it. With the new reinforcement signal, each unit gets a different local reward that evaluates its own action instead of the entire team's action, thus allowing efficient structural credit assignment. This idea of measuring the worth of a unit by its outgoing weight's norm has a long history (Uhr and Vossler 1961; Klopf and Gose 1969; Selfridge 1988).

With the new reinforcement signal, the learning rule at time t becomes:

$$\mathbf{w}^{(l)} \leftarrow \mathbf{w}^{(l)} + \Delta \mathbf{w}_t^{(l)}, \quad (7)$$

$$\Delta \mathbf{w}_t^{(l)} = \alpha G_t^{(l)} \nabla_{\mathbf{w}^{(l)}} \log \pi_l(H_t^{(0:l-1)}, H_t^{(l)}; \mathbf{w}^{(l)}), \quad (8)$$

where $1 \leq l \leq L$. Note that the only difference of the above learning rule with (3) is the change of the reinforcement signals to hidden units. We call this new learning method *Weight Maximization*. To apply the learning rule, we compute $\Delta \mathbf{w}_t^{(l)}$ iteratively from $l = L$ to 1. The pseudo-code can be found in Algorithm 1 of Appendix B. The computational cost of Weight Maximization is the same as backpropagation since it is linear in the number of layers.

To illustrate the algorithm, we consider an example of a simple network with only two units (one hidden unit and one output unit), as shown in Figure 1. For every time step t , the algorithm performs the following steps:

1. Given state S_t , we sample the activation value of hidden unit $H_t^{(1)} \sim \pi_1(S_t, \cdot; w^{(1)})$. For example, if the unit is a Bernoulli-logistic unit and the state S_t is one-dimensional, then $\Pr(H_t^{(1)} = 1 | S_t) = \sigma(w^{(1)} S_t)$ and $\Pr(H_t^{(1)} = 0 | S_t) = 1 - \sigma(w^{(1)} S_t)$, where σ is the sigmoid function. We sample output unit $H_t^{(2)} \sim \pi_2(H_t^{(1)}, \cdot; w^{(2)})$ similarly, and pass action $A_t = H_t^{(2)}$ to the environment.
2. After receiving return G_t from the environment (which is only known after the end of episode, but can be replaced with TD error for online learning), we use it to train the output unit by REINFORCE: $w^{(2)} = w^{(2)} + \alpha G_t^{(2)} \nabla_{w^{(2)}} \log \pi_2(H_t^{(1)}, H_t^{(2)}; w^{(2)})$, where α is the step size and $G_t^{(2)} = G_t$.
3. We compute the reinforcement signal to the hidden unit by $G_t^{(1)} = v^{(1)} \Delta v^{(1)} = w^{(2)}(w^{(2)} - w^{(2)})$, which is used to train the hidden unit by REINFORCE: $w^{(1)} = w^{(1)} + \alpha G_t^{(1)} \nabla_{w^{(1)}} \log \pi_1(S_t, H_t^{(1)}; w^{(1)})$.

In the following, we discuss the theoretical properties of Weight Maximization.

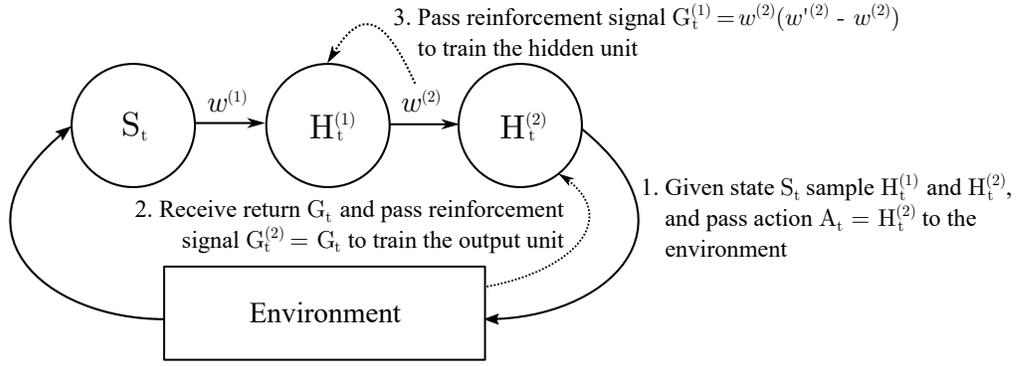


Figure 1: Illustration of a network with two units trained by Weight Maximization. See text for explanation.

Goal Alignment Condition

In this section, we address an important question: Under which situations is the goal of the hidden units and the whole network aligned? To simplify the discussion, we only consider single-time-step MDPs in this section and drop the subscript t , but the theorems here can be generalized to multiple-time-step MDPs.

To understand when the hidden units are maximizing the global reward, we analyze the gradient followed by the learning rule of the hidden units. First, the output unit is maximizing the return G as a result of Theorem 1 and REINFORCE:

$$\mathbb{E}[\Delta \mathbf{w}^{(L)}] \propto \nabla_{\mathbf{w}^{(L)}} \mathbb{E}[G]. \quad (9)$$

Then consider the learning rule of unit $L-1$, which is maximizing $G^{(L-1)}$:

$$\mathbb{E}[\Delta \mathbf{w}^{(L-1)}] \propto \nabla_{\mathbf{w}^{(L-1)}} \mathbb{E}[G^{(L-1)}] \quad (10)$$

$$= \nabla_{\mathbf{w}^{(L-1)}} \mathbb{E}[\mathbf{v}^{(L-1)} \cdot \Delta \mathbf{v}^{(L-1)}]. \quad (11)$$

$$\propto \nabla_{\mathbf{w}^{(L-1)}} (\mathbf{v}^{(L-1)} \cdot \nabla_{\mathbf{v}^{(L-1)}} \mathbb{E}[G]). \quad (12)$$

The last line is due to the fact that $\Delta \mathbf{v}^{(L-1)}$ is an entry in the vector $\Delta \mathbf{w}^{(L)}$, and so we can substitute the expectation of it with (9). We can continue the same process to derive the formulas of $\mathbb{E}[\Delta \mathbf{w}^{(l)}]$ for $l = L-2, L-3, \dots, 1$. This shows that the learning rule of hidden units is related to high-order cross partial derivatives of $\mathbb{E}[G]$ instead of the first-order derivative. To have the goal of units and the network aligned, it is sufficient and necessary that the cross partial derivatives are the same as the first-order derivative. Formally,

Lemma 1. *Let the policy be a network of stochastic units as defined above, and $\Delta \mathbf{w}^{(l)}$ be defined by (8). Then $\mathbb{E}[\Delta \mathbf{w}^{(l)}] \propto \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G]$ for all $1 \leq l \leq L-1$ if and only if $\nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \mathbb{E}[G]) \propto \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G]$ for all $1 \leq l \leq L-1$.*

The proof can be found in Appendix A.1. Therefore, we define the *goal alignment condition*, which is a sufficient condition that the hidden units are also maximizing the global return when applying Weight Maximization, by:

Definition 1. *Let the policy be a network of stochastic units as defined above. We say that the network has satisfied the goal alignment condition in an MDP if for all $1 \leq l \leq L-1$,*

$$\nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \mathbb{E}[G]) \propto \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G].$$

Except in special cases like a piecewise linear network with a piecewise linear reward function, this goal alignment condition does not hold exactly. However, the goal alignment condition holds approximately for all networks without extra assumptions:

Theorem 2. *Let the policy be a network of stochastic units as defined above. For all $1 \leq l \leq L-1$,*

$$\nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \mathbb{E}[G]) = \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G] + \mathcal{O}(\|\mathbf{v}^{(l)}\|_2^2).$$

The proof can be found in Appendix A.2. Therefore, the learning rule of hidden unit l is only approximately following the gradient of return in expectation with an error of $\sum_{k=l+1}^L \mathcal{O}(\|\mathbf{w}^{(k)}\|_2^2)$, since the error accumulates across units. In other words, the bias associated with the learning rule scales with the L^2 norm of the unit's outgoing weight. It is interesting to see that the more 'successful' a unit is (measured by the norm of its outgoing weight), the more severe is the problem of goal misalignment.

To combat the problem of goal misalignment, we suggest adding L^2 regularization, or weight decay, which can be seen as a soft constraint on the L^2 norm of weights (Goodfellow, Bengio, and Courville 2016) and thus prevents weights from having a large magnitude.

By replacing the rewards to the hidden units with (4) and adding weight regularization, the original cooperative game is turned into a *competitive game*. A competitive game refers to the scenario where each agent is receiving different rewards (Sutton and Barto 2018). With weight regularization, the upper layer has a 'limited' norm of weight to allocate due to the soft constraint, and so the units in the lower layer have to compete for the limited resources. In other words, units want to maximize their outgoing weights' norm, but the outgoing units want to minimize their incoming weights' norm, and therefore competition exists between units.

Weight Maximization with Eligibility Traces

Though Weight Maximization does not require a separate feedback pathway, the learning of a hidden unit needs to wait for the outgoing weight to finish updating; in contrast, in biological neurons, the change of synaptic strength has a slower timescale than the activation of neurons (Nicoll 2017), making it difficult to be used as an immediate feedback signal.

	Multiplexer		CartPole		Acrobot		LunarLander	
	Mean	Std.	Mean	Std.	Mean	Std.	Mean	Std.
Weight Max	0.81	0.01	390.38	43.25	-97.05	2.90	111.23	16.58
Weight Max w/ traces	n.a.	n.a.	373.11	17.86	-105.00	4.38	39.04	14.39
REINFORCE	0.29	0.01	163.92	64.43	-134.15	8.87	-94.23	19.03
REINFORCE with Thomas (2011)	0.28	0.02	363.71	24.07	-112.26	8.85	-66.19	67.14
STE Backprop	0.76	0.01	420.73	18.30	-98.55	5.05	47.53	35.04
Backprop	0.84	0.01	411.82	25.98	-91.82	5.24	71.77	21.02

Table 1: Average return over all episodes.

Besides different timescales, the change in synaptic strength of a biological neuron is the result of the neuron’s activity within a time interval instead of a single discrete time step. For example, multiple pairs of spikes are required to induce noticeable change in synaptic weights in spike-timing-dependent plasticity (STDP) experiments (Citri and Malenka 2008; Gerstner et al. 2014). In short, to be closer to biologically-observed plasticity rules, the effect of a unit’s action on the outgoing weight should be slow and long-lasting, instead of immediate and precise as assumed in Weight Maximization. We propose to use eligibility traces to solve both issues.

In the following discussion, we only consider a multi-layer network of stochastic units consisting of $M \geq 1$ layers. The L units in the network are arranged into a multi-layer structure such that the weights connecting units on non-adjacent layers are frozen to zero. We also denote $d(l)$ as the layer in which unit l resides and assume the last layer only contains the output unit.

Assume each unit requires a single time step to compute the weight update. That is, the weight update of unit l at time t , denoted by $\Delta \mathbf{w}_t^{(l)}$, is based on $G_{t-1}^{(l)}$ (defined in (4)) and $H_{t-1}^{(l)}$, the reward and action at the previous time step. Since $\Delta \mathbf{w}_t^{(l)}$ does not depend on the weight change at the same time step, we can compute it for all layers in parallel. However, the reward for the hidden layer m is lagging behind by $M - m$ time steps, as it takes a single time step for each of the $M - m$ upper layers to compute their weight updates. Therefore, the first issue can be seen as the problem of delayed reward - the action of a hidden unit on layer m at time t affects its rewards at time $t + M - m$, but not before.

For the second issue, consider the scenario where the effect of a unit’s action is long-lasting on the change of its outgoing weight. For example, if the output unit learns with decaying eligibility traces (Sutton and Barto 2018), then all $\Delta \mathbf{w}_{t+1}^{(L)}, \Delta \mathbf{w}_{t+2}^{(L)}, \dots$ depends on the value of layer $M - 1$ at time t , though the dependence decays with time. In other words, the action of layer $M - 1$ at time t will affect the change in its outgoing weight at and after time $t + 1$. We can continue the discussion for layer $M - 2, M - 3, \dots, 1$. From this perspective, the second issue can again be seen as the problem of delayed reward - the action of a hidden unit on layer m at time t affects all its rewards at and after time $t + M - m$, but not before.

The problem of delayed reward is well studied in RL, and one prominent and one biologically plausible solution is eligibility traces. We suggest using the following decay function $\lambda^l(t) : \mathbb{Z} \rightarrow [0, 1]$ for the unit l :

$$\lambda^l(t) = \begin{cases} 0 & \text{for } t \leq M - d(l) - 1, \\ (1 - \lambda)\lambda^{t - (M - d(l))} & \text{else,} \end{cases} \quad (13)$$

where $\lambda \in [0, 1]$ is the decay rate. Therefore, $\lambda^l(t)$ is the exponentially decaying trace but shifted by $M - d(l)$ time steps, since the action of unit l does not affect the reward in the next $M - d(l)$ time steps. With this decay function, we can generalize Weight Maximization to use eligibility traces. The learning rule of *Weight Maximization with eligibility traces* for a multi-layer network of stochastic units at time t is given by :

$$\mathbf{w}^{(l)} \leftarrow \mathbf{w}^{(l)} + \Delta \mathbf{w}_t^{(l)}, \quad (14)$$

$$\Delta \mathbf{w}_t^{(l)} = \alpha \delta_t^{(l)} \mathbf{z}_{t-1}^{(l)}, \quad (15)$$

$$\delta_t^{(l)} = \begin{cases} \mathbf{v}^{(l)} \cdot \Delta \mathbf{v}_t^{(l)} & \text{for } l \leq L - 1, \\ R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t) & \text{for } l = L, \end{cases} \quad (16)$$

$$\mathbf{z}_t^{(l)} = \sum_{k=0}^{t-1} \lambda^l(k) \gamma^{k - (M - d(l))} \nabla_{\mathbf{w}^{(l)}} \log \pi_l(H_{t-k}^{(0:l-1)}, H_{t-k}^{(l)}; \mathbf{w}^{(l)}), \quad (17)$$

where α denotes the step size and $1 \leq l \leq L$. We also replace the return G_t by TD error $R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t)$ to make the algorithm online. In practice, the state values $V^\pi(S_t)$ is generally unknown, but we can estimate the state values by another network implementing a TD algorithm, which is called a critic network (Sutton and Barto 2018). The pseudo-code can be found in Algorithm 2 of Appendix B.

With the above modification, the effect of a unit’s action on its outgoing weight is slow and long-lasting, in accordance with biologically-observed plasticity rules. Weight Maximization with eligibility traces also solves the three problems of backprop discussed in the introduction. Its learning rule is local and does not require any separate feedback pathways. The algorithm can be implemented in parallel for all layers, and there are no distinct feedforward and feedback phases for the whole network.

It should be noted that the feedback signal in backprop can also be computed based on the change of a unit’s outgoing weight in some cases, eliminating the need for a separate feedback pathway. However, as discussed earlier, it is not biologically plausible to use the change of a unit’s outgoing weight as an immediate feedback signal. The solution presented in this section, namely eligibility traces, can only be applied to Weight Maximization but not backprop, since backprop requires the activation values to be precisely matched with the feedback signals at the same time step. This underlines one major difference between the two algorithms.

Related Work

Research on solving tasks by a team of RL agents has a long history. Narendra and Thathachar (1974, 2012) described the stochastic learning automata, and Klopf (1972, 1982) proposed the hedonistic neuron hypothesis, which conjectured that individual neurons seek to maximize their own pleasure, and the collective behavior of these neurons can yield powerful adaptive systems. Barto and Anandan (1985) and Barto (1985) introduced the $A_{R-\lambda P}$ algorithm and showed that a team of $A_{R-\lambda P}$ units could learn with a globally-broadcast reward signal. Extending this class of learning rule, Williams (1992) introduced REINFORCE, a special case of $A_{R-\lambda P}$ when $\lambda = 0$, and proved that a team of agents trained with REINFORCE ascends the average reward gradient. Such a team of agents is recently called coagent networks (Thomas 2011). Theories relating to training coagent networks have been investigated (Thomas 2011; Kostas, Nota, and Thomas 2020; Thomas and Barto 2011), and Thomas (2011) proposed a variance reduction method for training coagent networks with REINFORCE by disabling exploration randomly. Chung (2021) proposed the MAP propagation algorithm, which minimizes the energy of the network before applying REINFORCE, to reduce the variance efficiently. However, in these papers, each agent in the team receives the same reward signal. In contrast, we propose that each agent maximizes the norm of its outgoing weight instead of the same reward signal, which turns the problem from a *cooperative game* into a *competitive game*. Zhang, Yang, and Başar (2021) reviewed recent development in the wider field of multi-agent RL.

Measuring the worth of a unit by the norm of its outgoing weight has been proposed (Uhr and Vossler 1961; Klopf and Gose 1969; Selfridge 1988). In these papers, hidden units with small outgoing weights are replaced by new hidden units with random incoming weights. These methods are thus based on the evolution approach instead of the RL approach as in Weight Maximization. Anderson (1986) proposed A_{R-P} algorithm with Penalty Prediction, which pushes units with small outgoing weights to match the incoming values. Bucket Brigade (Holland 1985) algorithm uses the strength of a classifier to do credit assignment, but with a winner-take-all selection and in classifier systems.

In addition, there is a large literature on methods for training a network of stochastic units, and a review can be found in Weber et al. (2019). STE backprop (Bengio, Léonard, and Courville 2013) is a practical method of training a network

of stochastic discrete units. Though STE backprop does not follow the gradient of the loss function, it is arguably the most effective way of training quantized ANN (Courbariaux, Bengio, and David 2015; Rastegari et al. 2016) and Yin et al. (2018) provides some theoretical justification for STE backprop. However, STE backprop suffers the same problem with backprop regarding biological plausibility.

Besides a team of agents trained by REINFORCE, many biologically plausible alternatives to backprop have been proposed. Biologically plausible learning rules based on reward prediction errors and attentional feedback have been proposed (Pozzi, Bohte, and Roelfsema 2020; Roelfsema and Ooyen 2005; Rombouts, Bohte, and Roelfsema 2015); but these learning rules mostly require a non-local feedback signal. Moreover, local learning rules based on contrastive divergence or nudging the values of output units have been proposed (Movellan 1991; Hinton 2002; Scellier and Bengio 2017). See Lillicrap et al. (2020) for a comprehensive review of algorithms that approximate backprop with local learning rules based on the differences in units’ values. Contrary to these papers, Weight Maximization does not require any separate feedback pathways or distinct phases in learning. Also, most of these algorithms are applied in supervised or unsupervised learning tasks, while Weight Maximization is applied in RL tasks.

Experiments

We applied our algorithms to four RL tasks: multiplexer, CartPole, Acrobot, and LunarLander. Details of the tasks can be found in Appendix C. We tested two variants of Weight Maximization: i. Weight Maximization, ii. Weight Maximization with Eligibility Traces. For both variants, we used Algorithm 2 in Appendix B (i. corresponds to the case $\lambda = 0$). We did not test Weight Maximization with Eligibility Traces on the multiplexer task since the task only has a single time step.

All networks considered have the same architecture: a three-layer network of stochastic units, with the first hidden layer having 64 units, the second hidden layer having 32 units, and the output layer being a softmax layer. All hidden units are Bernoulli-logistic units, i.e. $\pi_l(h^{(0:l-1)}, h^{(l)}, \mathbf{w}^{(l)}) = \sigma(\mathbf{w}^{(l)} \cdot h^{(0:l-1)})^{h^{(l)}} (1 - \sigma(\mathbf{w}^{(l)} \cdot h^{(0:l-1)}))^{1-h^{(l)}}$, where σ is the sigmoid function.

We consider four baselines to train hidden units: i. REINFORCE, ii. REINFORCE with the variance reduction method proposed by Thomas (2011), iii. STE backprop (Bengio, Léonard, and Courville 2013) and iv. backprop. For i. to iii., the networks are the same as the one used in Weight Maximization. For iv., the Bernoulli-logistic units are replaced with deterministic rectified linear units (ReLUs) so the hidden units can be trained by backprop. In all baselines, the output unit was trained by REINFORCE, and we used eligibility traces. Thus, all learning methods in the experiments are variants of Actor-Critic with Eligibility Traces (episodic) (Sutton and Barto 2018) with different methods of accumulating trace (for the baselines) or different reward signals to each unit (for Weight Maximization). For the critic networks in all experiments, we used a three-

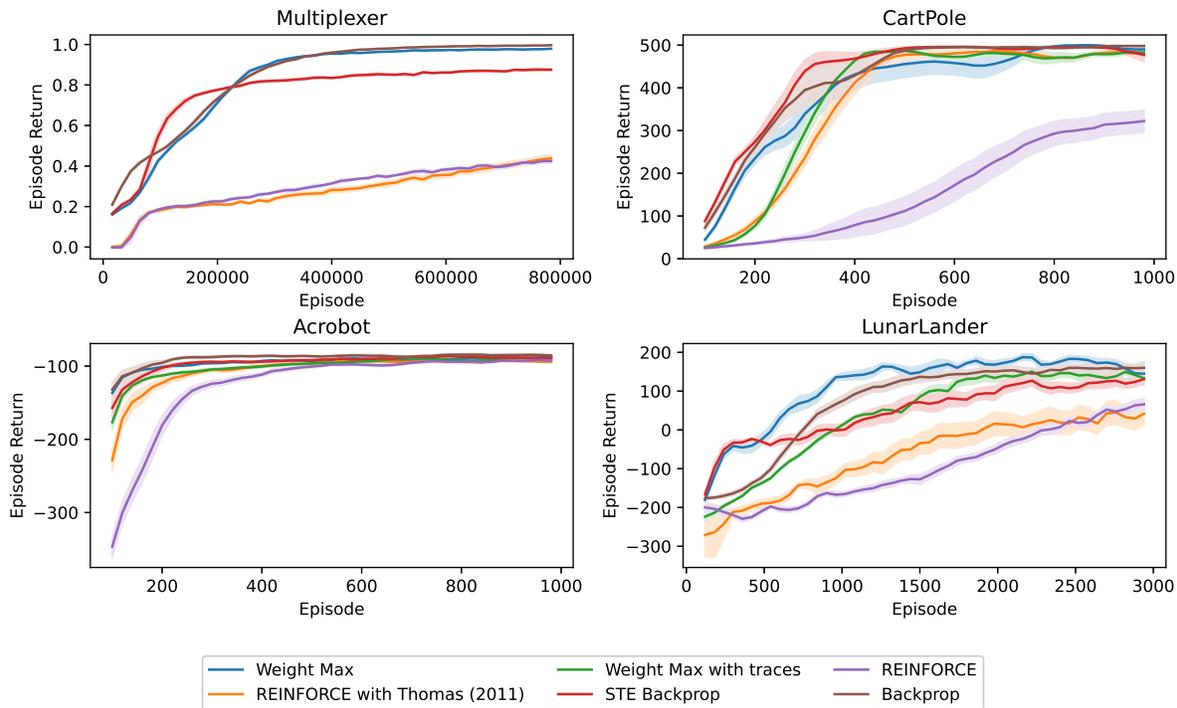


Figure 2: Episode returns in different RL tasks. Results are averaged over 10 independent runs, and shaded areas represent standard deviation over the runs. Curves are smoothed with a running average of 100 episodes (10,000 episodes for the multiplexer task).

layer ANN trained by backprop. Other experiments’ details and choice of hyperparameters can be found in Appendix C. The code is available at https://github.com/stephen-chung-mh/weight_max.

It should be noted that there are other algorithms besides actor-critic networks that can solve the tasks more efficiently, and the three-layer architecture is also chosen prior to experiments, which is not the optimal architecture since a shallower network performs better for simple tasks. The experiments only aim to compare different training methods for hidden units in a multi-layer actor network.

The average return over ten independent runs is shown in Fig 2. The mean and standard deviation of the average return can be found in Table 1. We observe that both variants of Weight Maximization have a significantly better performance than REINFORCE, suggesting that Weight Maximization allows effective structural credit assignment. Also, Weight Maximization has a similar performance compared to STE backprop, indicating that it is an effective method for training discrete units.

However, compared to a network of continuous-valued units trained by backprop, Weight Maximization performed slightly worse (except for the LunarLander). This is likely due to the limitation of discrete-valued units - units can only communicate with binary values instead of real values. This view is supported by the observation that discrete-valued units trained by STE backprop also performed worse than backprop. However, discrete units have the advantages

of low memory and communication costs; thus the slower learning may be compensated by a more efficient computation.

We notice that adding traces to Weight Maximization does not improve the performance. This is likely due to the harder temporal credit assignment. By using traces, the change in a unit’s outgoing weight at time t is affected by the unit’s action at time $t, t-1, t-2, \dots$ instead of only time t , making temporal credit assignment more difficult and thus learning slower.

In addition, we found that the representation learned by Weight Maximization is more statistically independent than backprop. For example, after training an agent to solve Cart-Pole by Weight Maximization (STE backprop), the absolute correlation across units on the first and the second layer is 0.690 (0.888) and 0.640 (0.830) respectively on average. This may be explained by the dynamics of Weight Maximization - units compete with each other to be more ‘useful’ and so have to learn different signals.

To better understand how different learning rules scale with the number of units in a network, we repeated the experiments on the multiplexer tasks with varying numbers of units in the network. Let the network have $64M$ and $32M$ units on the first and the second layer of the network respectively. The experimental results with different M can be found in Figure 3 and Figure 4.

We observe that the performance of both Weight Maximization and backprop improves with a larger M , while the

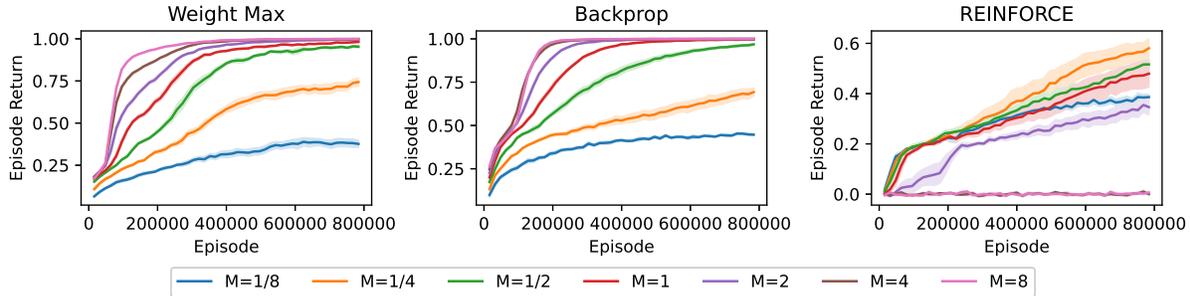


Figure 3: Episode returns in the multiplexer task with a varying number of units in the network. Results are averaged over 10 independent runs, and shaded areas represent standard deviation over the runs. Curves are smoothed with a running average of 10,000 episodes.

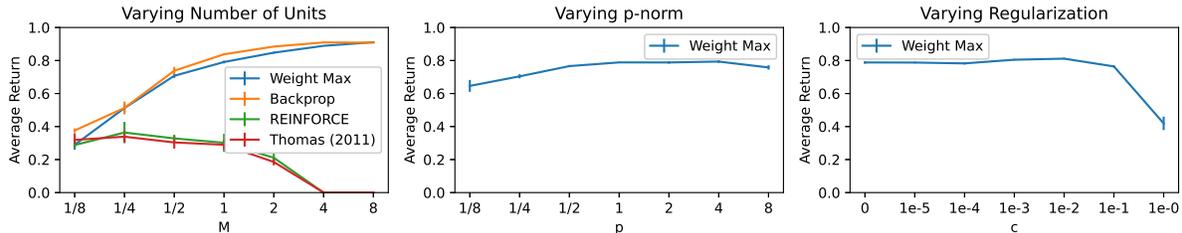


Figure 4: Average returns of all episodes in the multiplexer task with different hyperparameters. Results are averaged over 10 independent runs.

performance of REINFORCE peaks at $M = 1/4$. This illustrates a critical issue of REINFORCE that renders it impractical to train large networks - as all units are independently exploring and a single reward signal is used to evaluate the collective exploration of all units, it is more difficult to assign credit to each unit with a larger network. Thus, a larger network leads to a higher variance of parameters' updates and slower learning. This issue is alleviated in Weight Maximization since each unit receives a different reward signal that is strongly correlated with the unit's own action.

In Weight Maximization, we replace the reward to each unit by the change in the L^2 norm of the outgoing weight as shown in (4). It is also possible to use L^p norm instead of L^2 norm by generalizing (4) (see details in Appendix C). The experimental results on the multiplexer tasks with varying p can be found in Figure 4. The results show that the network is also able to learn with a similar performance when using other L^p norms.

We tested the effects of L^2 weight regularization, with $c \geq 0$ being the strength of weight regularization. That is, we add $-2cw^{(l)}$ to the learning rule (7). The experimental results on the multiplexer tasks with varying c can be found in Figure 4. We observe that the performance peaks at $c = 0.01$, suggesting that a small weight regularization can improve performance. This is in line with our analysis of the goal alignment condition.

We note that it is inefficient to apply Weight Maximization directly to train continuous-valued units. Additional mechanisms, such as equipping each unit with a baseline,

are required to match the speed of backprop and are left for future work.

Future Work and Conclusion

The approximate equivalence of every unit maximizing the global return in a cooperative game and every unit maximizing the norm of its outgoing weight in a competitive game offers a wide range of possible methods to train a multi-layer network efficiently besides backprop. Since each hidden unit faces a local optimization problem from the alternative perspective, learning rules besides REINFORCE can be applied to train each hidden unit. The lack of central coordination in Weight Maximization also leads to the possibility of implementing the algorithm asynchronously and efficiently with neuromorphic circuits (Indiveri et al. 2011).

In conclusion, we propose a novel algorithm that reduces the variance associated with training a team of agents with REINFORCE and thus significantly increases the learning speed. The proposed algorithm solves several major problems of backprop relating to biological plausibility. We also analyze the theoretical properties of the proposed algorithm and establish the approximate equivalence of hidden units maximizing the outgoing weights and the external rewards.

Acknowledgment

We would like to thank Andrew G. Barto, who inspired this research and provided valuable insights and comments.

References

- Anderson, C. W. 1986. *Learning and problem solving with multilayer connectionist systems*. Ph.D. thesis, Citeseer.
- Barto, A. G. 1985. Learning by statistical cooperation of self-interested neuron-like computing elements. *Human Neurobiology*, 4(4): 229–256.
- Barto, A. G.; and Anandan, P. 1985. Pattern-recognizing stochastic learning automata. *IEEE Transactions on Systems, Man, and Cybernetics*, (3): 360–375.
- Bengio, Y.; Lee, D.-H.; Bornschein, J.; Mesnard, T.; and Lin, Z. 2015. Towards biologically plausible deep learning. *arXiv preprint arXiv:1502.04156*.
- Bengio, Y.; Léonard, N.; and Courville, A. 2013. Estimating or propagating gradients through stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*.
- Brockman, G.; Cheung, V.; Pettersson, L.; Schneider, J.; Schulman, J.; Tang, J.; and Zaremba, W. 2016. Openai gym. *arXiv preprint arXiv:1606.01540*.
- Chung, S. 2021. MAP Propagation Algorithm: Faster Learning with a Team of Reinforcement Learning Agents. *Advances in Neural Information Processing Systems*, 34.
- Citri, A.; and Malenka, R. C. 2008. Synaptic plasticity: multiple forms, functions, and mechanisms. *Neuropsychopharmacology*, 33(1): 18–41.
- Courbariaux, M.; Bengio, Y.; and David, J.-P. 2015. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in neural information processing systems*, 3123–3131.
- Crick, F. 1989. The recent excitement about neural networks. *Nature*, 337(6203): 129–132.
- Gerstner, W.; Kistler, W. M.; Naud, R.; and Paninski, L. 2014. *Neuronal dynamics: From single neurons to networks and models of cognition*. Cambridge University Press.
- Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. *Deep learning*. MIT press.
- Hassabis, D.; Kumaran, D.; Summerfield, C.; and Botvinick, M. 2017. Neuroscience-inspired artificial intelligence. *Neuron*, 95(2): 245–258.
- Hinton, G. E. 2002. Training products of experts by minimizing contrastive divergence. *Neural computation*, 14(8): 1771–1800.
- Holland, J. H. 1985. Properties of the bucket brigade. In *Proceedings of the 1st International Conference on Genetic Algorithms*, 1–7.
- Indiveri, G.; Linares-Barranco, B.; Hamilton, T. J.; Van Schaik, A.; Etienne-Cummings, R.; Delbruck, T.; Liu, S.-C.; Dudek, P.; Häfliger, P.; Renaud, S.; et al. 2011. Neuro-morphic silicon neuron circuits. *Frontiers in neuroscience*, 5: 73.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Klopf, A. H. 1972. *Brain function and adaptive systems: a heterostatic theory*. 133. Air Force Cambridge Research Laboratories, Air Force Systems Command, United States Air Force.
- Klopf, A. H. 1982. *The hedonistic neuron: a theory of memory, learning, and intelligence*. Toxicology-Sci.
- Klopf, A. H.; and Gose, E. 1969. An evolutionary pattern recognition network. *IEEE Transactions on Systems Science and Cybernetics*, 5(3): 247–250.
- Kostas, J.; Nota, C.; and Thomas, P. 2020. Asynchronous Coagent Networks. In *International Conference on Machine Learning*, 5426–5435. PMLR.
- Lillicrap, T. P.; Santoro, A.; Marris, L.; Akerman, C. J.; and Hinton, G. 2020. Backpropagation and the brain. *Nature Reviews Neuroscience*, 21(6): 335–346.
- Mazzoni, P.; Andersen, R. A.; and Jordan, M. I. 1991. A more biologically plausible learning rule for neural networks. *Proceedings of the National Academy of Sciences*, 88(10): 4433–4437.
- Movellan, J. R. 1991. Contrastive Hebbian learning in the continuous Hopfield model. In *Connectionist models*, 10–17. Elsevier.
- Narendra, K. S.; and Thathachar, M. A. 1974. Learning automata—a survey. *IEEE Transactions on systems, man, and cybernetics*, (4): 323–334.
- Narendra, K. S.; and Thathachar, M. A. 2012. *Learning automata: an introduction*. Courier corporation.
- Nicoll, R. A. 2017. A brief history of long-term potentiation. *Neuron*, 93(2): 281–290.
- O’Reilly, R. C. 1996. Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural computation*, 8(5): 895–938.
- Pozzi, I.; Bohte, S.; and Roelfsema, P. 2020. Attention-Gated Brain Propagation: How the brain can implement reward-based error backpropagation. *Advances in Neural Information Processing Systems*, 33.
- Rastegari, M.; Ordonez, V.; Redmon, J.; and Farhadi, A. 2016. Xnor-net: Imagenet classification using binary convolutional neural networks. In *European conference on computer vision*, 525–542. Springer.
- Roelfsema, P. R.; and Ooyen, A. v. 2005. Attention-gated reinforcement learning of internal representations for classification. *Neural computation*, 17(10): 2176–2214.
- Rombouts, J. O.; Bohte, S. M.; and Roelfsema, P. R. 2015. How attention can create synaptic tags for the learning of working memories in sequential tasks. *PLoS Comput Biol*, 11(3): e1004060.
- Rumelhart, D. E.; Hinton, G. E.; and Williams, R. J. 1986. Learning representations by back-propagating errors. *Nature*, 323(6088): 533–536.
- Scellier, B.; and Bengio, Y. 2017. Equilibrium propagation: Bridging the gap between energy-based models and backpropagation. *Frontiers in computational neuroscience*, 11: 24.
- Selfridge, O. G. 1988. Pandemonium: A paradigm for learning. In *Neurocomputing: Foundations of research*, 115–122. A Bradford Book.
- Sutton, R. S.; and Barto, A. G. 2018. *Reinforcement learning: An introduction*. MIT press.

- Thomas, P. S. 2011. Policy gradient coagent networks. In *Advances in Neural Information Processing Systems*, 1944–1952.
- Thomas, P. S.; and Barto, A. G. 2011. Conjugate Markov Decision Processes. In *International Conference on Machine Learning*, 137–144.
- Uhr, L.; and Vossler, C. 1961. A pattern recognition program that generates, evaluates, and adjusts its own operators. In *Papers presented at the May 9-11, 1961, western joint IRE-AIEE-ACM computer conference*, 555–569.
- Weber, T.; Heess, N.; Buesing, L.; and Silver, D. 2019. Credit assignment techniques in stochastic computation graphs. In *The 22nd International Conference on Artificial Intelligence and Statistics*, 2650–2660. PMLR.
- Williams, R. J. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4): 229–256.
- Yin, P.; Lyu, J.; Zhang, S.; Osher, S.; Qi, Y.; and Xin, J. 2018. Understanding Straight-Through Estimator in Training Activation Quantized Neural Nets. In *International Conference on Learning Representations*.
- Zhang, K.; Yang, Z.; and Başar, T. 2021. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of Reinforcement Learning and Control*, 321–384.

A. Proof

A.1 Proof of Lemma 1

Proof. Assume $\mathbb{E}[\Delta \mathbf{w}^{(l)}] \propto \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G]$ for all $1 \leq l \leq L-1$. Let $1 \leq l \leq L-1$. Then,

$$\nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G] \tag{18}$$

$$\propto \mathbb{E}[\Delta \mathbf{w}^{(l)}] \tag{19}$$

$$\propto \mathbb{E}[(\mathbf{v}^{(l)} \cdot \Delta \mathbf{v}^{(l)}) \nabla_{\mathbf{w}^{(l)}} \log \pi_l(H^{(0:l-1)}, H^{(l)}; \mathbf{w}^{(l)})] \tag{20}$$

$$= \nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \mathbb{E}[\Delta \mathbf{v}^{(l)}]) \tag{21}$$

$$\propto \nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \mathbb{E}[G]). \tag{22}$$

Conversely, assume $\nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \mathbb{E}[G]) \propto \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G]$ for all $1 \leq l \leq L-1$. We prove by induction on k : For $k = L$, $\mathbb{E}[\Delta \mathbf{w}^{(k)}] \propto \nabla_{\mathbf{w}^{(k)}} \mathbb{E}[G]$. Let $1 \leq k < L$ and assume $\mathbb{E}[\Delta \mathbf{w}^{(l)}] \propto \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G]$ for all $k+1 \leq l \leq L$. Then,

$$\nabla_{\mathbf{w}^{(k)}} \mathbb{E}[G] \tag{23}$$

$$\propto \nabla_{\mathbf{w}^{(k)}} (\mathbf{v}^{(k)} \cdot \nabla_{\mathbf{v}^{(k)}} \mathbb{E}[G]) \tag{24}$$

$$\propto \nabla_{\mathbf{w}^{(k)}} (\mathbf{v}^{(k)} \cdot \mathbb{E}[\Delta \mathbf{v}^{(k)}]) \tag{25}$$

$$= \mathbb{E}[(\mathbf{v}^{(k)} \cdot \Delta \mathbf{v}^{(k)}) \nabla_{\mathbf{w}^{(k)}} \log \pi_l(H^{(0:k-1)}, H^{(k)}; \mathbf{w}^{(k)})] \tag{26}$$

$$\propto \mathbb{E}[\Delta \mathbf{w}^{(k)}]. \tag{27}$$

□

A.2 Proof of Theorem 2

Proof. Let $1 \leq l \leq L-1$. Then,

$$\nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \mathbb{E}[G]) \tag{28}$$

$$= \nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) \mathbb{E}[G | H^{(0:l)} = h^{(0:l)}]) \tag{29}$$

$$= \nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \nabla_{\mathbf{v}^{(l)}} \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)})) \tag{30}$$

$$= \nabla_{\mathbf{w}^{(l)}} (\mathbf{v}^{(l)} \cdot \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) h^{(l)} \nabla g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)})) \tag{31}$$

$$= \nabla_{\mathbf{w}^{(l)}} \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) (h^{(l)} \mathbf{v}^{(l)} \cdot \nabla g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)})) \tag{32}$$

$$= \nabla_{\mathbf{w}^{(l)}} \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) (g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)}) - g(\mathbf{0}; h^{(0:l-1)}) + \mathcal{O}(\|\mathbf{v}^{(l)}\|_2^2)) \tag{33}$$

$$= \nabla_{\mathbf{w}^{(l)}} \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)}) + \mathcal{O}(\|\mathbf{v}^{(l)}\|_2^2) \tag{34}$$

$$= \nabla_{\mathbf{w}^{(l)}} \sum_{h^{(0:l)}} \Pr(H^{(0:l)} = h^{(0:l)}) \mathbb{E}[G | H^{(0:l)} = h^{(0:l)}] + \mathcal{O}(\|\mathbf{v}^{(l)}\|_2^2) \tag{35}$$

$$= \nabla_{\mathbf{w}^{(l)}} \mathbb{E}[G] + \mathcal{O}(\|\mathbf{v}^{(l)}\|_2^2). \tag{36}$$

(29) to (30) uses the fact that $\mathbb{E}[G | H^{(0:l)} = h^{(0:l)}]$ can be expressed as a differentiable function $g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)})$:

$$\mathbb{E}[G | H^{(0:l)} = h^{(0:l)}] \tag{37}$$

$$= \sum_{h^{(l+1:L)}} \mathbb{E}[G | S = h^{(0)}, A = h^{(L)}] \pi_L(h^{(0:L-1)}; h^{(L)}; \mathbf{w}^{(L)}) \pi_{L-1}(h^{(0:L-2)}, h^{(L-1)}; \mathbf{w}^{(L-1)}) \dots \pi_{l+1}(h^{(0:l)}, h^{(l+1)}; \mathbf{w}^{(l+1)}) \tag{38}$$

$$= \sum_{h^{(l+1:L)}} \mathbb{E}[G | S = h^{(0)}, A = h^{(L)}] f_L(\mathbf{w}^{(L)} \cdot h^{(0:L-1)}, h^{(L)}) f_{L-1}(\mathbf{w}^{(L-1)} \cdot h^{(0:L-2)}, h^{(L-1)}) \dots f_{l+1}(\mathbf{w}^{(l+1)} \cdot h^{(0:l)}, h^{(l+1)}) \tag{39}$$

$$= g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)}), \tag{40}$$

where (38) to (39) uses the assumption that for all l , $\pi_l(h^{(0:l-1)}, h^{(l)}; \mathbf{w}^{(l)})$ can be expressed as a function $f_l(\mathbf{w}^{(l)} \cdot h^{(0:l-1)}, h^{(l)})$.

(32) to (33) uses the Taylor's formula: $g(\mathbf{0}; h^{(0:l)}) = g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)}) - h^{(l)} \mathbf{v}^{(l)} \cdot \nabla g(h^{(l)} \mathbf{v}^{(l)}; h^{(0:l-1)}) + \mathcal{O}(\|\mathbf{v}^{(l)}\|_2^2)$.

(33) to (34) uses the fact that $g(\mathbf{0}; h^{(0:l-1)})$ does not depend on $h^{(l)}$ thus the derivative of its expectation w.r.t. $\mathbf{w}^{(l)}$ is 0. □

Algorithm 1: Weight Maximization

```
1 Input: differentiable policy function:  $\pi_l(h^{(0:l-1)}, h^{(l)}; \mathbf{w}^{(l)})$  for  $l \in \{1, 2, \dots, L\}$ ;  
2 Algorithm Parameter: step size  $\alpha > 0$ ; discount rate  $\gamma \in [0, 1]$ ;  
3 Initialize policy parameter:  $\mathbf{w}^{(l)}$  for  $l \in \{1, 2, \dots, L\}$ ;  
4 Loop forever (for each episode):  
5   Generate an episode  $S_0, H_0, R_1, \dots, S_{T-1}, H_{T-1}, R_T$  following  $\pi_l(\cdot, \cdot; \mathbf{w}^{(l)})$  for  $l \in \{1, 2, \dots, L\}$ ;  
6   Loop for each step of the episode,  $t = 0, 1, \dots, T - 1$ :  
7      $G^{(L)} \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$ ;  
8     for  $l = L, L - 1, \dots, 1$  do  
9        $\Delta \mathbf{w}^{(l)} \leftarrow G^{(l)} \nabla_{\mathbf{w}^{(l)}} \log \pi_l(H_t^{(0:l-1)}, H_t^{(l)}; \mathbf{w}^{(l)})$ ;  
10       $\mathbf{w}^{(l)} \leftarrow \mathbf{w}^{(l)} + \alpha \Delta \mathbf{w}^{(l)}$ ;  
11       $G^{(l-1)} \leftarrow \sum_{k=l}^L \mathbf{w}^{(k)}_{(l-1)} \Delta \mathbf{w}^{(k)}_{(l-1)}$ ;  
12    end
```

Algorithm 2: Weight Maximization with Eligibility Traces

```
1 Input: differentiable policy function:  $\pi_m(h^{(m-1)}, h^{(m)}; W^{(m)})$  for  $m \in \{1, 2, \dots, M\}$ ;  
2 Algorithm Parameter: step size  $\alpha > 0$ ; trace decay rate  $\lambda \in [0, 1]$ ; discount rate  $\gamma \in [0, 1]$ ;  
3 Initialize policy parameter:  $W^{(m)}$  for  $l \in \{1, 2, \dots, M\}$ ;  
4 Loop forever (for each episode):  
5   Initialize  $S_1$  (first state of episode) ;  
6   Initialize zero eligibility traces  $\mathbf{z}^{(m)}$  and zero  $\Delta W^{(m)}$  for  $m \in \{1, 2, \dots, M\}$  ;  
7   Loop while  $S_t$  is not terminal (for each time step  $t = 1, 2, \dots$ ):  
8      $H_t^0 \leftarrow S_t$  ;  
9     /* Sample Action */  
10    Sample  $H_t^{(m)}$  from  $\pi_m(H_t^{(m-1)}, \cdot; W^{(m)})$  for  $m \in \{1, 2, \dots, M\}$  ;  
11    /* Compute  $\delta^{(m)}$  and apply REINFORCE */  
12    if episode not in first time step then  
13      Receive TD error  $\delta$  from the critic network;  
14       $\delta^{(M)} \leftarrow \delta$ ;  
15       $\delta^{(m)} \leftarrow (W^{(m+1)} \odot \Delta W^{(m+1)})^T \mathbf{1}$  for  $m \in \{1, 2, \dots, M - 1\}$ ;  
16       $\Delta W^{(m)} \leftarrow (\delta^{(m)} \mathbf{1}^T) \odot \mathbf{z}^{(m)}$  for  $m \in \{1, 2, \dots, M\}$ ;  
17       $W^{(m)} \leftarrow W^{(m)} + \alpha \Delta W^{(m)}$  for  $m \in \{1, 2, \dots, M\}$ ;  
18    end  
19    /* Trace accumulation */  
20     $\mathbf{z}^{(m)} \leftarrow \gamma \lambda \mathbf{z}^{(m)} + \nabla_{W^{(m)}} \log \pi_m(H_{t-M+m}^{(m-1)}, H_{t-M+m}^{(m)}; W^{(m)})$  for  $m \in \{1, 2, \dots, M\}$ ;  
21    Take action  $H_t^{(M)}$ , observe  $S_{t+1}, R_{t+1}$  ;
```

Table 2: The hyperparameters of Weight Maximization used in the experiments.

	Multiplexer	CartPole	Acrobot	LunarLander
α_1	1e-2	2e-2	1e-1	8e-2
α_2	1e-3	2e-4	1e-3	8e-4
α_3	1e-4	2e-5	1e-4	8e-5
T	1	1	1	2
γ	n.a.	.98	.98	.98

B. Algorithms

The notations in Algorithm 2 are slightly different from the main paper: for $1 \leq m \leq M$, we let $H_t^{(m)}$ denote the values of units on layer m at time t , which is a multivariate random variable. The distribution of $H_t^{(m)}$ conditioned on $H_t^{(m-1)}$ is given by $\Pr(H^{(m)} = h^{(m)} | H^{(m-1)} = h^{(m-1)}) = \pi_m(h^{(m-1)}, h^{(m)}; W^{(m)})$, where $W^{(m)} \in \mathbb{R}^{n(m) \times n(m-1)}$ is the incoming weights for layer m and $n(m)$ denotes the number of units on layer m . With this new formulation, (16) can be expressed as:

$$\delta^{(m)} = \begin{cases} (W^{(m+1)} \odot \Delta W_t^{(m+1)})^T \mathbf{1} & \text{for } m \in \{1, 2, \dots, M-1\}, \\ R_t + \gamma V^\pi(S_{t+1}) - V^\pi(S_t) & \text{for } m = M, \end{cases}$$

where \odot is elementwise multiplication and $\mathbf{1}$ is a vector of ones. Note that $\delta^{(m)} \in \mathbb{R}^{n(m)}$ for $m < M$ instead of \mathbb{R} as in original notation, since there are $n(m)$ different rewards to each $n(m)$ units on layer m . Also, (15) can be expressed as:

$$\Delta W_t^{(m)} = (\delta_{t-1}^{(m)} \mathbf{1}^T) \odot \mathbf{z}_{t-1}^{(m)}$$

which is equivalent to multiplying each row i of $\mathbf{z}_{t-1}^{(m)}$ by the entry i of $\delta_{t-1}^{(m)}$.

Note that in line 17 of Algorithm 2, the trace accumulation is delayed by $M - m$ time steps, and so the last $M - m$ actions of layer m have to be stored in a buffer. If $t - M + m < 1$, then the gradient is replaced with 0.

C. Implementation Details of Experiments

In the multiplexer task, the state is sampled from all possible values of a binary vector of size $k + 2^k$ with equal probability. The action set is $\{-1, 1\}$, and we give a reward of +1 if the action of the agent is the desired action and -1 otherwise. The desired action is given by the output of a k -bit multiplexer, with the input of the k -bit multiplexer being the state. We consider $k = 4$ here, so the dimension of the state space is 20. This is similar to the 2-bit multiplexer considered by Barto (1985). For the implementation of the remaining three tasks, we used CartPole-v1, Acrobot-v1, and LunarLander-v2 in OpenAI’s Gym (Brockman et al. 2016).

The hyperparameters of Weight Maximization used in the experiments can be found in Table 2. We used a different learning rate α for each layer of the network, and we denote the learning rate for the l^{th} layer to be α_l . For the learning rule in line 15 of the pseudo-code, we used Adam optimizer (Kingma and Ba 2014) instead, with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We did not use any batch update. We denote the temperature of the output softmax layer to be T . These hyperparameters were selected based on manual tuning to optimize the learning curve. We did the same manual tuning for the hyperparameters of the baseline models. For softmax layers, the last logit unit is set to be a constant zero since $n - 1$ logits are sufficient to represent a distribution of a random variable with n possible values.

For all learning methods with eligibility traces (including the baseline models), we used $\lambda = 0.8$ except for the LunarLander task, where we used $\lambda = 0.9$. The value of λ was selected prior to the experiments, similar to the choice of the network architecture.

We did not use critic networks for all experiments in the multiplexer task and used the reward from the environment as the reinforcement signal since there is only a single time step in the task. For the critic networks in other tasks, we used a three-layer ANN, with the same architecture as the actor network but with ReLu hidden units and a linear output unit. The critic network was trained by backprop as in Actor-Critic with Eligibility Traces (episodic) (Sutton and Barto 2018).

We annealed the learning rate linearly such that the learning rate is reduced to 1e-2 and 1e-1 of the initial learning rates at 5e4 and 1e6 steps in CartPole and Acrobot respectively, and the learning rate remains unchanged afterward. We also annealed the learning rate linearly for the baseline models. We found that this can make the final performance more stable. We did not anneal the learning rate in the other two tasks.

We use the following formula to generalize Weight Maximization to use L^p norm instead of L^2 norm (note that this generalizes (4)):

$$G_t^{(l)} = \begin{cases} (\nabla_{\mathbf{v}^{(l)}} \|\mathbf{v}^{(l)}\|_p^p) \cdot \Delta \mathbf{v}_t^{(l)} & \text{for } l \in \{1, 2, \dots, L-1\}, \\ G_t & \text{for } l = L. \end{cases} \quad (41)$$