
RUDDER: Return Decomposition for Delayed Rewards

Jose A. Arjona-Medina* Michael Gillhofer* Michael Widrich*
Thomas Unterthiner Sepp Hochreiter
LIT AI Lab
Institute for Machine Learning
Johannes Kepler University Linz, Austria

Abstract

1 We propose a novel reinforcement learning approach for finite Markov decision
2 processes (MDPs) with delayed rewards. In this work, biases of temporal difference
3 (TD) estimates are proved to be corrected only exponentially slowly in the number
4 of delay steps. Furthermore, variances of Monte Carlo (MC) estimates are proved
5 to increase the variance of other estimates, which number can exponentially grow
6 in the number of delay steps. We introduce RUDDER, a return decomposition
7 method, which creates a new MDP with same optimal policies as the original
8 MDP but with redistributed rewards that have largely reduced delays. If the return
9 decomposition is optimal, then the new MDP does not have delayed rewards and
10 TD estimates are unbiased. In this case, the rewards track Q -values so that the
11 future expected reward is always zero. We experimentally confirm our theoretical
12 results on bias and variance of TD and MC estimates. On artificial tasks with
13 different lengths of reward delays, we show that RUDDER is exponentially faster
14 than TD, MC, and MC Tree Search (MCTS). RUDDER outperforms rainbow, A3C,
15 DDQN, Distributional DQN, Dueling DDQN, Noisy DQN, and Prioritized DDQN
16 on the delayed reward Atari game Venture in only a fraction of the learning time.
17 RUDDER considerably improves the state-of-the-art on the delayed reward Atari
18 game Bowling in much less learning time.

19 1 Introduction

20 Assigning the credit for a received reward to actions that were performed, is one of the central tasks
21 in reinforcement learning [58]. Long term credit assignment has been identified as one of the largest
22 challenges in reinforcement learning [46]. Current reinforcement learning methods are still slowed
23 down significantly when facing long-delayed rewards [41, 30]. To learn delayed rewards there are
24 three phases to consider: (1) discovering the delayed reward, (2) keeping information about the
25 delayed reward, (3) learning to receive the delayed reward to secure it for the future. Recent successful
26 reinforcement methods provide solutions to one or more of these items. Most prominent are Deep
27 Q -Networks (DQNs) [32, 33], which combine Q -learning with convolutional neural networks for
28 visual reinforcement learning [24]. The success of DQNs is attributed to *experience replay* [29],
29 which stores observed state-reward transitions and then samples from them. Prioritized experience
30 replay [47, 22] advanced the sampling from the replay memory. Different policies perform exploration
31 in parallel for the Ape-X DQN and share a prioritized experience replay memory [22]. IMPALA
32 improves A2C by parallel actors and corrects for policy-lags between actors and learners [10]. DQN
33 was extended to double DQN (DDQN) [60, 61] which helps exploration as the overestimation bias
34 is reduced. Noisy DQNs [11] explore by a stochastic layer in the policy network (see [18, 48]).

*authors contributed equally

35 Distributional Q -learning [6] profits from noise since means that have high variance are more likely
 36 selected. The dueling network architecture [62, 63] separately estimates state values and action
 37 advantages, which helps exploration in unexperienced states. Policy gradient approaches [66] like
 38 A3C with asynchronous gradient descent [31] or Ape-X DPG [22] explore via parallel policies,
 39 too. Proximal policy optimization (PPO) extends A3C by a surrogate objective and a trust region
 40 optimization realized by clipping or a Kullback-Leibler distance penalty [50].

41 Recent approaches aim to solve learning problems caused by delayed rewards. Function approxima-
 42 tions of value functions or critics [33, 31] bridge time intervals if states associated with rewards are
 43 similar to states that were encountered many steps earlier. For example, assume a function that learned
 44 to predict a large reward at the end of an episode if a state has a particular feature. The function
 45 can generalize this correlation to the begin of an episode and predict already high reward for states
 46 possessing the same feature. Multi-step temporal difference (TD) learning [56, 58] improved both
 47 DQNs and policy gradients [17, 31]. AlphaGo and AlphaZero learned to play Go and Chess better
 48 than human professionals using Monte Carlo Tree Search (MCTS) [51, 52]. MCTS simulates games
 49 from a time point until the end of the game or an evaluation point, therefore captures long-delayed
 50 rewards. Recently, world models using a evolution strategy were successful [14]. These forward
 51 view approaches using world models are not feasible in probabilistic environments with a high state
 52 transition branching factor. Backward view approaches trace back from known goal states [9] or from
 53 high-reward states [13]. However a step-by-step backward model has to be learned.

54 We propose learning from a backward view, which is constructed from a forward model. The forward
 55 model predicts the return, while the backward analysis identifies states and actions which have
 56 caused the return. We apply Long Short-Term Memory (LSTM) [19, 21] to predict the return of
 57 an episode. LSTM was already used in reinforcement learning [49] for advantage learning [3] and
 58 learning policies [15, 31, 16]. However sensitivity analysis by “backpropagation through a model”
 59 [35, 44, 45, 4] has major drawbacks: local minima, instabilities, exploding or vanishing gradients in
 60 the world model, proper exploration, contribution (relevance) of actions are not regarded only their
 61 sensitivity [18, 48].

62 Since sensitivity analysis substantially hinders learning, we use contribution analysis for backward
 63 analysis like contribution-propagation [25], contribution approach [38], excitation backprop [68],
 64 layer-wise relevance propagation (LRP) [2], Taylor decomposition [2, 34], or integrated gradients
 65 (IG) [55]. Using contribution analysis, a predicted return can be decomposed into contributions along
 66 the state-action sequence. Substituting the prediction by the actual return, we obtain a redistributed
 67 reward leading to new MDP with the same optimal policies as for the original MDP. Redistributing
 68 the reward is fundamentally different from reward shaping [36, 64], which changes the reward
 69 as a function of states but not of actions. Reward redistribution is related to “look-back advice”
 70 [65] which, in contrast to reward redistribution, still requires the original MDP for learning. We
 71 propose RUDDER, which performs reward redistribution by return decomposition and, therefore,
 72 overcomes problems of TD and MC stemming from delayed rewards. RUDDER vastly decreases
 73 the variance of MC and largely avoids the exponentially slow bias corrections of TD — for optimal
 74 return decomposition TD is even unbiased.

75 2 Bias-Variance for MDP Estimates

76 We perform a bias-variance analysis for temporal difference (TD) and Monte Carlo (MC) estimators of
 77 the action-value function. A finite Markov decision process (MDP) \mathcal{P} is 6-tuple $\mathcal{P} = (\mathcal{S}, \mathcal{A}, \mathcal{R}, p, \pi, \gamma)$
 78 of finite sets \mathcal{S} of states s (random variable S_t at time t), \mathcal{A} of actions a (random variable A_t),
 79 and \mathcal{R} of rewards r (random variable R_{t+1}). Furthermore, \mathcal{P} has transition-reward distributions
 80 $p(S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a)$ conditioned on state-actions, a policy given as an action
 81 distributions $\pi(A_{t+1} = a' \mid S_{t+1} = s')$ conditioned on states, and a discount factor $\gamma \in [0, 1]$.
 82 The marginals are $p(r \mid s, a) = \sum_{s'} p(s', r \mid s, a)$ and $p(s' \mid s, a) = \sum_r p(s', r \mid s, a)$. The
 83 expected reward is $r(s, a) = \sum_r r p(r \mid s, a)$. The return G_t is $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$. We often
 84 consider finite horizon MDPs with sequence length T and $\gamma = 1$ giving $G_t = \sum_{k=0}^{T-t} R_{t+k+1}$. The
 85 action-value function $q^\pi(s, a)$ for policy π is $q^\pi(s, a) = \mathbb{E}_\pi[G_t \mid S_t = s, A_t = a]$. Goal of learning
 86 is to maximize the expected return at time $t = 0$, that is $v_0^\pi = \mathbb{E}_\pi[G_0]$.

87 **Bias-Variance Analysis for MDP Estimates.** MC estimates $q^\pi(s, a)$ by an arithmetic mean of the
 88 return, while TD methods like SARSA or Q -learning estimate $q^\pi(s, a)$ by an exponential average of
 89 the return. When using Monte Carlo for learning a policy we use an exponential average, too, since the
 90 policy steadily changes. The i th update of action-value q at state-action (s_t, a_t) is $(q^\pi)^{i+1}(s_t, a_t) =$
 91 $(q^\pi)^i(s_t, a_t) + \alpha \left(\sum_t^T r_{t+1} - (q^\pi)^i(s_t, a_t) \right)$. Assume n samples $\{X_1, \dots, X_n\}$ from a distribution
 92 with mean μ and variance σ^2 . For these samples, we compute bias and variance of the arithmetic
 93 mean $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$ and the exponential average $\tilde{\mu}_n = \alpha \sum_{i=1}^n (1 - \alpha)^{n-i} X_i + (1 - \alpha)^n \mu_0$
 94 with μ_0 as initial value and $\alpha \in (0, 1)$. We obtain $\text{bias}(\hat{\mu}_n) = 0$ and $\text{var}(\hat{\mu}_n) = \sigma^2/n$ as well
 95 as $\text{bias}(\tilde{\mu}_n) = (1 - \alpha)^n (\mu_0 - \mu)$ and $\text{var}(\tilde{\mu}_n) = \sigma^2 (\alpha(1 - (1 - \alpha)^{2n})) / (2 - \alpha)$ (see Appendix
 96 A1.2.1 for more details). Both variances are proportional to σ^2 , which is the variance when sampling
 97 a return from the MDP \mathcal{P} .

98 Using $\mathbb{E}_{s', a'}(f(s', a')) = \sum_{s'} p(s' | s, a) \sum_{a'} \pi(a' | s') f(s', a')$, and analog $\text{Var}_{s', a'}$ and Var_r , the
 99 next theorem gives mean and variance $V^\pi(s, a) = \text{Var}_\pi[G_t | s, a]$ of sampling returns from an MDP.

100 **Theorem 1.** *The mean q^π and variance V^π of sampled returns from an MDP are*

$$q^\pi(s, a) = \sum_{s', r} p(s', r | s, a) \left(r + \gamma \sum_{a'} \pi(a' | s') q^\pi(s', a') \right) = r(s, a) + \gamma \mathbb{E}_{s', a'} [q^\pi(s', a') | s, a],$$

$$V^\pi(s, a) = \text{Var}_r[r | s, a] + \gamma^2 (\mathbb{E}_{s', a'} [V^\pi(s', a') | s, a] + \text{Var}_{s', a'} [q^\pi(s', a') | s, a]). \quad (1)$$

101 The proof is given after Theorem A1 in the appendix. The theorem extends the deterministic
 102 reward case [54, 59]. The variance $V^\pi(s, a)$ consists of three parts: (i) The immediate vari-
 103 ance $\text{Var}_r[r | s, a]$ stemming from the probabilistic reward $p(r | s, a)$. (ii) The local variance
 104 $\gamma^2 \text{Var}_{s', a'} [q^\pi(s', a') | s, a]$ caused by probabilistic state transitions and probabilistic policy. (iii)
 105 The expected variance $\gamma^2 \mathbb{E}_{s', a'} [V^\pi(s', a') | s, a]$ of the next Q -values, which is zero for TD since
 106 it replaces $q^\pi(s', a')$ by fixed $\hat{q}^\pi(s', a')$. Therefore TD has less variance than MC which uses the
 107 complete future return. See Appendix A1.2.2 for more details.

108 **Delayed Reward Aggravates Learning.** The i th temporal difference update with learning rate α
 109 of the action-value $q(s_t, a_t)$ is

$$q^{i+1}(s_t, a_t) = q^i(s_t, a_t) + \alpha (r_{t+1} + \mathcal{A}_{a'}(q^i(s_{t+1}, a')) - q^i(s_t, a_t)), \quad (2)$$

110 with $\mathcal{A}_{a'}(\cdot) = \max_{a'}(\cdot)$ (Q -learning), $\mathcal{A}_{a'}(\cdot) = \sum_{a'} \pi(a' | s_{t+1})(\cdot)$ (expected SARSA), $\mathcal{A}_{a'}(\cdot)$
 111 sample a' from $\pi(a' | s_{t+1})$ (SARSA). The next theorem states that TD has an exponential decay for
 112 Q -value updates even for eligibility traces [23, 5, 57, 53].

113 **Theorem 2.** *For initialization $q^0(s_t, a_t) = 0$ and delayed reward with $r_t = 0$ for $t \leq T$,
 114 $q(s_{T-i}, a_{T-i})$ receives its first update not earlier than at episode i via $q^i(s_{T-i}, a_{T-i}) = \alpha^{i+1} r_{T+1}^1$,
 115 where r_{T+1}^1 is the reward of episode 1. Eligibility traces with $\lambda \in [0, 1)$ lead to an exponential decay
 116 of $(\gamma\lambda)^k$ when the reward is propagated k steps back.*

117 The proof is given after Theorem A2 in the appendix. To correct the TD bias by a certain amount
 118 requires exponentially many updates with the number of delay steps.

119 For Monte Carlo the variance of a single delayed reward can increase the variance of action-values of
 120 all previously visited state-actions. We define the “on-site” variance ω

$$\omega(s, a) = \text{Var}_r[r | s, a] + \text{Var}_{s', a'} [q^\pi(s', a') | s, a], \quad (3)$$

121 \mathbf{V}^π is the vector with value $V^\pi(s', a')$ at position (s', a') and \mathbf{P}_t the transition matrix from states
 122 s_t to s_{t+1} with entries $p(s_{t+1} | s_t, a_t) \pi(a_{t+1} | s_{t+1})$ at position $((s_t, a_t), (s_{t+1}, a_{t+1}))$. For finite
 123 time horizon, the “backward induction algorithm” [39, 40] gives with $\mathbf{V}_{T+1}^\pi = \mathbf{0}$ and $\omega_{T+1} = \mathbf{0}$ and
 124 row-stochastic matrix, $\mathbf{P}_{t \leftarrow k} = \prod_{\tau=t}^{k-1} \mathbf{P}_\tau$:

$$\mathbf{V}_t^\pi = \sum_{k=t}^T \prod_{\tau=t}^{k-1} \mathbf{P}_\tau \omega_k = \sum_{k=t}^T \mathbf{P}_{t \leftarrow k} \omega_k, \quad (4)$$

125 where we define $\prod_{\tau=t}^{t-1} \mathbf{P}_\tau = \mathbf{I}$ and $[\omega_k]_{(s_k, a_k)} = \omega(s_k, a_k)$. We are interested in the number of
 126 action-values which variances are affected through the increase of the variance of a single delayed

127 reward. Let N_t be the number of all states s_t that are reachable after t time steps of an episode. Let
 128 c_t be the random average connectivity of a state in s_t to states in s_{t-1} . Let n_t be number of states
 129 in s_t that are affected by ω_k for $t \leq k$ with $n_k = 1$ (only one action-value with delayed reward at
 130 time $t = k$). Next theorem says that the on-site variance ω_k can have large effects on the variance of
 131 action-values of all previously visited state-actions, which number can grow exponentially.

132 **Theorem 3.** For $t \leq k$, on-site variance ω_k at step k contributes to V_t^π by the term $\mathbf{P}_{t \leftarrow k} \omega_k$, where
 133 $\|\mathbf{P}_{t \leftarrow k}\|_\infty = 1$. The number a_k of states affected by ω_k is $a_k = \sum_{t=0}^k \left(1 - \left(1 - \frac{c_t}{N_{t-1}}\right)^{n_t}\right) N_{t-1}$.

134 The proof can be found after Theorem A3. For small k , the number a_k of states affected by on-site
 135 variance ω_k at step k grows exponentially with k . For large k and after some time $t > \hat{t}$, the number
 136 a_k of states affected by ω_k grows linearly. (See Corollary A1 in the appendix). Consequently, we
 137 aim for decreasing the on-site variance ω_k for large k , in order to reduce the variance. In summary,
 138 delayed rewards lead to exponentially slow corrections of biases of temporal difference (TD) and
 139 can increase exponentially many variances of Monte Carlo (MC) action-value estimates, where the
 140 exponentially grows is in both cases in the number of delay steps.

141 3 Return Decomposition and Reward Redistribution

142 A Markov decision process (MDP) $\tilde{\mathcal{P}}$ is *state-enriched* compared to a MDP \mathcal{P} if $\tilde{\mathcal{P}}$ has the same
 143 states, actions, transition probabilities, and reward probabilities as \mathcal{P} but with additional information
 144 in their states. We observe that \mathcal{P} is a homomorphic image of $\tilde{\mathcal{P}}$ with the same actions. Therefore
 145 each optimal policy $\tilde{\pi}^*$ of $\tilde{\mathcal{P}}$ has an equivalent optimal policy π^* of \mathcal{P} , and vice versa, with the same
 146 optimal return [42, 43]. These properties are known from state abstraction and aggregation [28] and
 147 from bisimulation [12]. For more details see Appendix A1.3.1. Two Markov decision processes $\tilde{\mathcal{P}}$
 148 and \mathcal{P} are *return-equivalent* if they differ only in $p(\tilde{r} | s, a)$ and $p(r | s, a)$ but for each policy π they
 149 have the same expected return at $t = 0$: $\tilde{v}_0^\pi = v_0^\pi$. Return-equivalent decision processes have the
 150 same optimal policies.

151 We assume to have an MDP \mathcal{P} with immediate reward which is transformed to a state-enriched
 152 MDP $\tilde{\mathcal{P}}$ with delayed reward, where the return is given as reward at sequence end. The transformed
 153 delayed state-enriched MDP has reward $\tilde{r}_t = 0$, $t \leq T$, and $\tilde{r}_{T+1} = \sum_{k=0}^T R_{k+1}$. The states are
 154 enriched by ρ which records the accumulated already received rewards, therefore $\tilde{s}_t = (s_t, \rho_t)$,
 155 where $\rho_t = \sum_{k=0}^{t-1} r_{k+1}$. We show in Proposition ?? that $\tilde{q}^{\tilde{\pi}}(\tilde{s}, a) = q^\pi(s, a) + \sum_{k=0}^{t-1} r_{k+1}$ for
 156 $\tilde{\pi}(a | \tilde{s}) = \pi(a | s)$. Thus, each immediate reward MDP can be transformed into a delayed reward
 157 MDP without changing the optimal policies.

158 Next we consider the opposite direction, where the delayed reward MDP $\tilde{\mathcal{P}}$ is given and we want
 159 to find an immediate reward MDP \mathcal{P} . \mathcal{P} should be return-equivalent to $\tilde{\mathcal{P}}$ and differ from $\tilde{\mathcal{P}}$ only
 160 by its reward distributions. We have to redistribute the final reward, which is the return, \tilde{r}_{T+1} to
 161 previous time steps, therefore we have to decompose the return into a sum of rewards at different
 162 time steps. To allow for a return decomposition, we predict the return \tilde{r}_{T+1} by a function g using the
 163 state-action sequence: $g((s, a)_{0:T}) = \tilde{r}_{T+1}$, where $(s, a)_{0:T}$ is the state-action sequence from $t = 0$
 164 to $t = T$. In a next step we decompose g into a sum: $g((s, a)_{0:T}) = \sum_{t=0}^T h(a_t, s_t)$, where h is the
 165 prediction contribution. Since $\tilde{\mathcal{P}}$ is an MDP, the reward can be predicted from (a_T, s_T) since s_T
 166 contains information about the already accumulated reward. Therefore we use a difference $\Delta(s, s')$
 167 between state s and its successor s' instead of s to avoid the Markov property in the input sequence.
 168 The difference Δ is assumed to make $(a_t, \Delta(s_t, s_{t+1}))$ statistically independent from each other in
 169 the sequence $(a, \Delta)_{0:T} = (a_0, \Delta(s_0, s_1), \dots, a_t, \Delta(s_t, s_{t+1}))$. The function g is decomposed by
 170 contribution analysis into a sum of h by $g((a, \Delta)_{0:T}) = \tilde{r}_{T+1} = \sum_{t=0}^T h(a_t, \Delta(s_t, s_{t+1}))$. The
 171 actual reward redistribution is $r_{t+1} = \tilde{r}_{T+1} h(a_t, \Delta(s_t, s_{t+1})) / g((a, \Delta)_{0:T})$ to ensure $\sum_{t=0}^T \tilde{r}_{t+1} =$
 172 $\tilde{r}_{T+1} = \sum_{t=0}^T r_{t+1}$.

173 If for partial sums $\sum_{\tau=0}^t h(a_\tau, \Delta(s_\tau, s_{\tau+1})) = \tilde{q}^\pi(s_t, a_t)$ holds, then the *return decomposition is*
 174 *optimal*. We have for $g((a, \Delta)_{0:T}) = \tilde{r}_{T+1}$ rewards $R_0 = h_0 = h(a_0, \Delta(s_0, s_1)) = \tilde{q}^\pi(s_0, a_0)$ and
 175 $R_t = h_t = h(a_t, \Delta(s_t, s_{t+1})) = \tilde{q}^\pi(s_t, a_t) - \tilde{q}^\pi(s_{t-1}, a_{t-1})$. The term $\tilde{q}^\pi(s_{t-1}, a_{t-1})$ introduces
 176 variance in R_t .

177 **Theorem 4.** *The MDP \mathcal{P} based on the redistributed reward given by an optimal return decomposition*
 178 *(I) has the same optimal policies as $\tilde{\mathcal{P}}$ of the delayed reward, and (II) the Q -values are given by*
 179 $q^\pi(s_t, a_t) = r(s_t, a_t) = \tilde{q}^\pi(s_t, a_t) - \mathbb{E}[\tilde{q}^\pi(s_{t-1}, a_{t-1}) \mid s_t, a_t]$.

180 The proof can be found after Theorem 4 in the appendix. In particular, when start-
 181 ing with zero initialized Q -values, then TD learning of \mathcal{P} is not biased at the be-
 182 ginning. For policy gradients with eligibility traces using $\lambda \in [0, 1]$ for G_t^λ
 183 [58], we have the expected updates $\mathbb{E}_\pi \left[\nabla_\theta \log \pi(a_t \mid s_t; \theta) \sum_{\tau=0}^{T-t} \lambda^\tau q^\pi(s_{t+\tau}, a_{t+\tau}) \right] =$
 184 $\mathbb{E}_\pi \left[\nabla_\theta \log \pi(a_t \mid s_t; \theta) \sum_{\tau=0}^{T-t} \lambda^\tau r(s_{t+\tau}, a_{t+\tau}) \right]$, where $r(s_t, a_t)$ is replaced during learning by
 185 a sample from R_t which is the redistributed reward for an episode.

186 **RUDDER: Return Decomposition using LSTM.** We introduce RUDDER “RetUrn Decomposi-
 187 tion for DELayed Rewards”, which performs return decomposition using a Long Short-Term Memory
 188 (LSTM) network for redistributing the original reward. RUDDER consists of (I) a safe exploration
 189 strategy, (II) a lessons replay buffer, and, most importantly, (III) an LSTM with contribution analysis
 190 for return decomposition. **(I) Safe exploration.** Exploration strategies should assure that LSTM
 191 receives training data with delayed rewards. Toward this end we introduce a new exploration strategy
 192 which initiates at a certain time in the episode an exploration sequence to discover delayed rewards.
 193 To avoid an early stop of the exploration sequence, we perform a safe exploration which avoids
 194 actions associated with low Q -values. Low Q -values hint at states with zero future reward where the
 195 agent gets stuck. Exploration parameters are starting time, length, and the action selection strategy
 196 with safety constraints. **(II) Lessons replay buffer.** If safe exploration discovers an episode with
 197 unexpected delayed reward, it is secured in a lessons replay buffer [29]. Episodic memory has been
 198 used for episodic control [27] and for episodic backward update to efficiently propagate delayed
 199 rewards [26]. Unexpected is indicated by a large prediction error of LSTM. Sampling from lessons
 200 replay buffer is done similar to prioritized experience replay. Episodes with larger error are more often
 201 sampled. **(III) LSTM and contribution analysis.** LSTM networks [19, 21]), are used to predict the
 202 return from an input sequence. LSTM solves the vanishing gradient problem [19, 20], which severely
 203 impedes credit assignment in recurrent neural networks, i.e. the correct identification of relevant but
 204 delayed input events. LSTM backward analysis is done through contribution analysis like layer-wise
 205 relevance propagation (LRP) [2], Taylor decomposition [2, 34], or integrated gradients (IG) [55].
 206 These methods identify the contributions of the inputs to the final prediction, therefore supply the
 207 return decomposition.

208 The LSTM return decomposition is optimal if LSTM predicts at every time step the expected
 209 final return. To push LSTM toward optimal return decomposition, we introduce continuous return
 210 predictions as auxiliary tasks, where the LSTM has to predict the final return during the sequence.
 211 Hyperparameters are when and how often LSTM predicts and how continuous prediction errors are
 212 weighted. Strictly monotonic LSTM architecture (see AppendixA4.3.1) can also ensure that LSTM
 213 decomposition is optimal.

214 4 Experiments

215 We use $\gamma = 1$ for delayed rewards in MDPs with finite time horizon or absorbing states which has
 216 been confirmed to be suited to long delays by meta-gradient reinforcement learning [67].

217 **Grid World:** RUDDER is tested on a grid world with delayed reward at the end of an episode. The
 218 MDP is a 7×7 grid with 3 special locations (*start*, *key* and *door*), and 4 actions (*up*, *down*, *left* and
 219 *right*). An episode begins with a random *start* and ends at the fixed *door* or after 25 time steps. *key*
 220 defines the minimal delay. If the agent visits *door* at time t , then the reward is $-t \cdot 0.1$ and increased
 221 by 10 if the agent has visited *key*, otherwise the reward is 0. To investigate how the delay affects bias
 222 and variance, Q -values are estimated by TD and MC for a random policy which assures that all states
 223 are visited. After computing the true Q -values by backward induction, we compare the bias, variance
 224 and mean squared error (MSE) of the estimators for the MDP with delayed reward and the new MDP
 225 obtained by RUDDER with optimal reward redistribution. Figure 1 shows that RUDDER has smaller
 226 number of Q -values with high variance than the original MDP, when learning MC estimators. It also
 227 shows that RUDDER corrects the bias faster than TD estimators in the original MDP. So far we kept
 228 the policy constant and focused on learning the action-value function. Next, we compare Q -learning,

229 Monte Carlo (MC), and Monte Carlo Tree Search (MCTS) at learning a policy for the grid world.
 230 Figure 2 shows the number of episodes required by different methods to learn a policy that achieves
 231 90% of the return of the optimal policy for different delays. Optimal reward redistribution speeds up
 learning a policy exponentially. More information is available in Appendix A5.1.1.

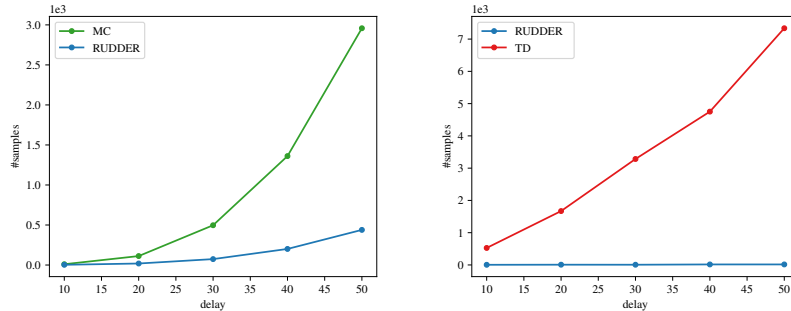


Figure 1: Experimental evaluation of MSE, bias, and variance of different Q -value estimators on the Grid World. Left: Variance of the MC estimator grows exponentially with the delay. Shown are the number of samples needed to go below a threshold. Right: Bias correction in TD is exponentially small with the delay. Shown are the number of samples needed to half the initial error for the initial state estimator in TD.

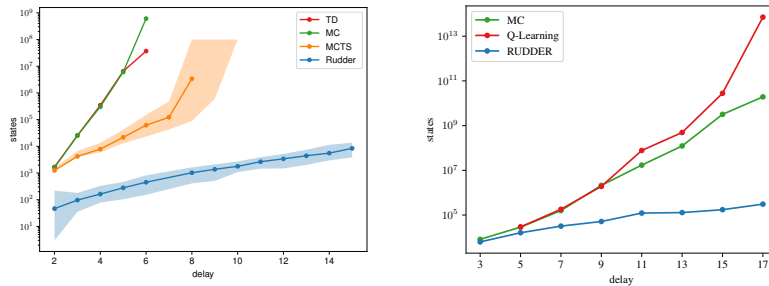


Figure 2: Number of observed states required by different methods to learn a policy that achieves 90% of the return of the optimal policy for different delays. We compare Q -learning, Monte Carlo (MC), and Monte Carlo Tree Search (MCTS). Left: Grid World environment. Right: Charge-Discharge environment. Reward redistribution requires an exponentially smaller number of states than the original methods.

232 **Charge-Discharge environment:** We test RUDDER on another task, the Charge-Discharge envi-
 234 ronment, which has two states: *discharged* D / *charged* C and two actions *discharge* d / *charge* c .
 235 The deterministic reward is $r(D, d) = 1$, $r(C, d) = 10$, $r(D, c) = 0$, and $r(C, c) = 0$. The reward
 236 $r(C, d)$ is accumulated for the whole episode and given only at time $T \in \{3, \dots, 13\}$, which deter-
 237 mines the maximal delay of a reward. The deterministic state transitions are $(\{D, C\}, d) \rightarrow D$ and
 238 $(\{D, C\}, c) \rightarrow C$. The optimal policy alternates between charging and discharging to accumulate
 239 a reward of 10 every other time step. RUDDER is based on a monotonic LSTM with layer-wise
 240 relevance propagation (LRP) for the backwards analysis (see Appendix ?? for more details). The
 241 reward redistribution provided by RUDDER served to learn a policy by Q -learning. We compare
 242 RUDDER with Q -learning, MC and MCTS. The results are shown in Figure 2. Reward redistribution
 243 requires an exponentially smaller number of states than Q -learning, MC and MCTS to learn the
 244 optimal policy.

245 **Atari Games Bowling and Venture:** We investigated the Atari games supported by the Arcade
 246 Learning Environment [7] and OpenAI Gym [8] for games with delayed reward. Requirements for
 247 proper games to demonstrate performance on delayed reward are: (I) large delay between an action
 248 and the resulting reward, (II) no distractions due to other rewards or changing characteristics of the
 249 environment, (III) no skills to be learned to receive the delayed reward. The requirements were met

250 by Bowling and Venture. In Bowling the only reward of the game is given at the end of the episode,
 251 200 frames after the first relevant action. In Venture the first reward has a minimum delay of 120
 252 frames from the first relevant action. Figure 3 shows that RUDDER learns faster than rainbow [17],
 253 Prioritized DDQN [47], Noisy DQN [11], Dueling DDQN [63], DQN [33], C51 (Distributional
 254 DQN) [6], DDQN [60], A3C [31], and Ape-X DQN [22]. RUDDER sets a new state-of-the-art score
 255 in Bowling after 12M environment frames. Thus, RUDDER outperforms its competitors in only 10%
 256 of their training time, as shown in Table 1. For more details see For more details see Appendix A5.2.

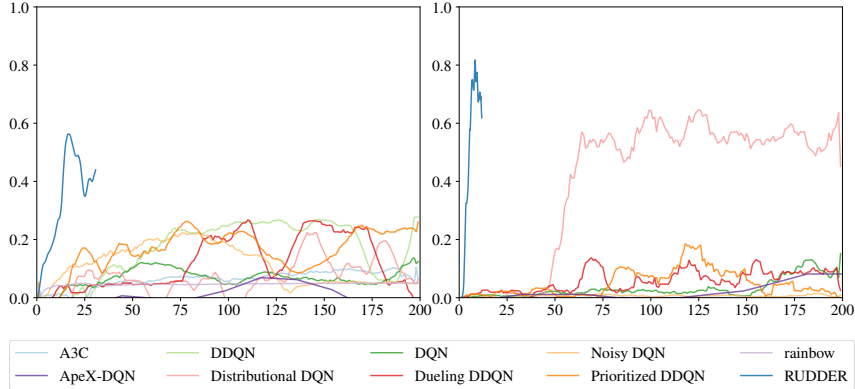


Figure 3: RUDDER learns the delayed reward for the Atari games Bowling and Venture faster than other methods. Normalized human-percentage scores during training for Bowling (left) and for Venture (right), where learning curves are taken from previous publications [17, 22]. RUDDER sets a new state-of-the-art for Bowling.

Algorithm	Frames	Bowling		Venture	
		%	raw	%	raw
RUDDER	12M	62.10	108.55	96.55	1,147
rainbow	200M	5.01	30	0.46	5.5
Prioritized DDQN	200M	28.71	62.6	72.67	863
Noisy DQN	200M	39.39	77.3	0	0
Dueling DDQN	200M	30.81	65.5	41.85	497
DQN	200M	19.84	50.4	13.73	163
Distributional DQN	200M	37.06	74.1	93.22	1,107
DDQN	200M	32.7	68.1	8.25	98
Ape-X DQN	22,800M	-17.6	4	152.67	1,813
Random	-	0	23.1	0	0
Human	-	100	160.7	100	1,187

Table 1: Results of RUDDER and other methods when learning the Atari games Bowling and Venture. Normalized human-percentage and raw scores over 200 testing-games with no-op starting condition: A3C scores are not reported, as not available for no-op starting condition. Scores for other methods were taken from previous publications [6, 17]. The RUDDER model is chosen based only on its training loss over 12M frames.

257 **RUDDER Implementation for Bowling and Venture.** We implemented RUDDER for the prox-
 258 imal policy optimization (PPO) algorithm [50]. For policy gradients the expected updates are
 259 $E_{\pi} [\nabla_{\theta} \log \pi(a | s; \theta) q^{\pi}(s, a)]$, where $q^{\pi}(s, a)$ is replaced during learning by the return G_t or its
 260 expectation. RUDDER policy gradients replace $q^{\pi}(s, a)$ by the redistributed reward $r(s, a)$ assuming
 261 an optimal return decomposition. With eligibility traces using $\lambda \in [0, 1]$ for G_t^{λ} [58], we have
 262 the rewards $\rho_t = r_t + \lambda \rho_{t+1}$ with $\rho_{T+1} = 0$ and the expected updates $E_{\pi} [\nabla_{\theta} \log \pi(a_t | s_t; \theta) \rho_t]$.
 263 We use integrated gradients [55] for the backward analysis of RUDDER. The LSTM prediction
 264 g is decomposed by the integrated gradient IG via the equation $g(x) = \text{IG}(g, x, \mathbf{0}) + g(\mathbf{0}) =$
 265 $\sum_{t=0}^T (h(a_t, \Delta(s_t, s_{t+1})) + (1/(T+1))g(\mathbf{0}))$. For Atari games, Δ is defined as pixel-wise differ-
 266 ence of two consecutive frames. To make static objects visible, we augment the input with the current
 267 frame. For more implementation details see Appendix A5.2. Source code will be made available.

268 **Evaluation Methodology.** Agents were trained for 12M environment frames with *no-op starting*
 269 *condition*, i.e. a random number of up to 30 no-operation actions at the start of a game. Training
 270 episodes are terminated by loss of life or at 108K frames. After training, the best model was
 271 selected based on training data and evaluated on 200 games with no-op starting condition and a
 272 maximum length of 108K frames, following [61]. For comparison across games, the normalized
 273 human-percentage scores according to [6] are reported.

274 **Visual Confirmation of the Learning Boost of Reward Redistribution.** We visually confirmed
 275 a meaningful and helpful redistribution of reward in both Bowling and Venture during training.
 276 As illustrated in Figure 4, RUDDER is capable of redistributing a reward to key events in game,
 277 drastically shortening the delay of the reward and quickly steering the agent toward good policies.
 Furthermore, it enriches sequences that were sparse in reward with a dense reward signal.

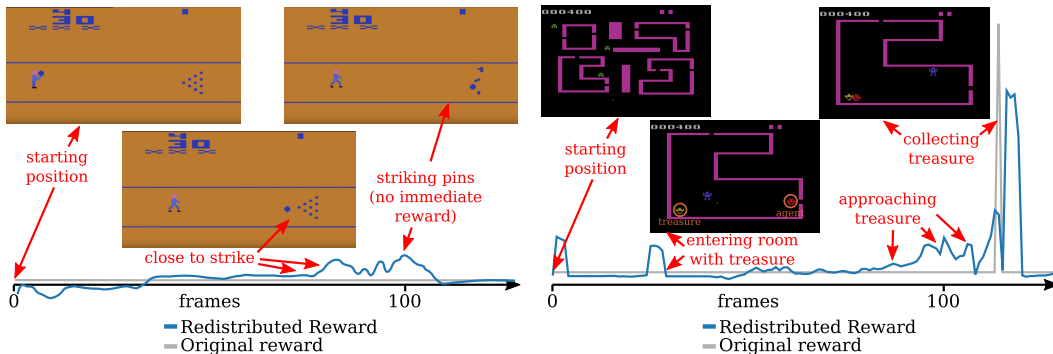


Figure 4: Observed return decomposition by RUDDER in two Atari games with long delayed rewards. **Left:** In the game Bowling reward is only given after three strikes have been performed. RUDDER identifies the actions that guide the ball in the right direction to hit all pins. Once the ball hit the pins, RUDDER detects the delayed reward associated with striking the pins down. In the figure only 100 frames are represented but the whole episode spans 200 frames. In the original game, the reward is given only at the end of the episode. **Right:** In the game Venture reward is only obtained after picking the treasure. RUDDER guides the agent (red) towards the treasure (golden) via reward redistribution. Reward is redistributed to entering a room with treasure. Furthermore, the redistributed reward gradually increases as the agent approaches the treasure. For illustration purposes, the green curve shows the return redistribution before applying lambda. The environment only gives reward at the event of collecting treasure (blue curve).

278

279 5 Discussion and Conclusion

280 **Exploration** is the most critical part of RUDDER, since discovering delayed rewards is the first step
 281 to exploit them.

282 **Human expert episodes** are an alternative to exploration and can serve to fill the lessons replay
 283 buffer. Learning can be sped up considerably when LSTM identifies human key actions. Return
 284 decomposition will reward human key actions even for episodes with low return since other actions
 285 that thwart high returns receive negative reward. Using human demonstrations in reinforcement
 286 learning led to a huge improvement on some Atari games like Montezuma’s Revenge [37, 1].

287 **Conclusion.** We have shown that for finite Markov decision processes with delayed rewards TD
 288 exponentially slowly corrects biases and MC can increase exponentially many variances of estimates,
 289 both in the number of delay steps. We have introduced RUDDER, a return decomposition method,
 290 which creates a new MDP that keeps the optimal policies but its redistributed rewards do not have
 291 delays. In the optimal case TD for the new MDP is unbiased. On two artificial tasks we demonstrated
 292 that RUDDER is exponentially faster than TD, MC, and MC Tree Search (MCTS). For the Atari
 293 game Venture with delayed reward RUDDER outperforms all methods except Ape-X DQN in
 294 much less learning time. For the Atari game Bowling with delayed reward RUDDER improves the
 295 state-of-the-art and outperforms PPO, Rainbow, and APE-X with less learning time.

296 **Acknowledgments**

297 This work was supported by NVIDIA Corporation, Bayer AG with Research Agreement 09/2017,
298 Merck KGaA, Zalando SE with Research Agreement 01/2016, Audi.JKU Deep Learning Center, Audi
299 Electronic Venture GmbH, Janssen Pharmaceutica, IWT research grant IWT150865 (Exaptation),
300 LIT grant LSTM4Drive, and FWF grant P 28660-N31.

301 **References**

- 302 [1] Y. Aytar, T. Pfaff, D. Budden, T. Le Paine, Z. Wang, and N. de Freitas. Playing hard exploration
303 games by watching YouTube. *ArXiv*, 2018.
- 304 [2] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek. On pixel-wise
305 explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLoS ONE*,
306 10(7):e0130140, 2015.
- 307 [3] B. Bakker. Reinforcement learning with long short-term memory. In T. G. Dietterich, S. Becker,
308 and Z. Ghahramani, editors, *Advances in Neural Information Processing Systems 14*, pages
309 1475–1482. MIT Press, 2002.
- 310 [4] B. Bakker. Reinforcement learning by backpropagation through an lstm model/critic. In *IEEE*
311 *International Symposium on Approximate Dynamic Programming and Reinforcement Learning*,
312 pages 127–134, 2007.
- 313 [5] A. G. Barto and R. S. Sutton. Landmark learning: An illustration of associative search.
314 *Biological Cybernetics*, 42(1):1–8, 1981.
- 315 [6] M. G. Bellemare, W. Dabney, and R. Munos. A distributional perspective on reinforcement
316 learning. In D. Precup and Y. W. Teh, editors, *Proceedings of the 34th International Conference*
317 *on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research (ICML)*, pages
318 449–458. PMLR, 2017.
- 319 [7] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. The Arcade learning environment: An
320 evaluation platform for general agents. *Journal of Artificial Intelligence Research*, 47:253–279,
321 2013.
- 322 [8] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
323 Openai gym. *ArXiv*, 2016.
- 324 [9] A. D. Edwards, L. Downs, and J. C. Davidson. Forward-backward reinforcement learning.
325 *ArXiv*, 2018.
- 326 [10] L. Espeholt, H. Soyer, R. Munos, K. Simonyan, V. Mnih, T. Ward, Y. Doron, V. Firoiu, T. Harley,
327 I. Dunning, S. Legg, and K. Kavukcuoglu. IMPALA: Scalable distributed Deep-RL with
328 importance weighted actor-learner architectures. In J. Dy and A. Krause, editors, *Proceedings*
329 *of the 35th International Conference on Machine Learning*, 2018. ArXiv: 1802.01561.
- 330 [11] M. Fortunato, M. G. Azar, B. Piot, J. Menick, I. Osband, A. Graves, V. Mnih, R. Munos,
331 D. Hassabis, O. Pietquin, C. Blundell, and S. Legg. Noisy networks for exploration. *ArXiv*,
332 2018. Sixth International Conference on Learning Representations (ICLR).
- 333 [12] R. Givan, T. Dean, and M. Greig. Equivalence notions and model minimization in Markov
334 decision processes. *Artificial Intelligence*, 147(1):163–223, 2003.
- 335 [13] A. Goyal, P. Brakel, W. Fedus, T. Lillicrap, S. Levine, H. Larochelle, and Y. Bengio. Recall
336 traces: Backtracking models for efficient reinforcement learning. *ArXiv*, 2018.
- 337 [14] D. Ha and J. Schmidhuber. World models. *ArXiv*, 2018.
- 338 [15] M. J. Hausknecht and P. Stone. Deep recurrent Q-Learning for partially observable MDPs.
339 *ArXiv*, 2015.
- 340 [16] N. Heess, G. Wayne, Y. Tassa, T. P. Lillicrap, M. A. Riedmiller, and D. Silver. Learning and
341 transfer of modulated locomotor controllers. *ArXiv*, 2016.

- 342 [17] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot,
343 M. G. Azar, and D. Silver. Rainbow: Combining improvements in deep reinforcement learning.
344 *ArXiv*, 2017.
- 345 [18] S. Hochreiter. Implementierung und Anwendung eines ‘neuronalen’ Echtzeit-Lernalgorithmus
346 für reaktive Umgebungen. Practical work, Supervisor: J. Schmidhuber, Institut für Informatik,
347 Technische Universität München, 1990.
- 348 [19] S. Hochreiter. Untersuchungen zu dynamischen neuronalen Netzen. Master’s thesis, Technische
349 Universität München, 1991.
- 350 [20] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber. Gradient flow in recurrent nets: the
351 difficulty of learning long-term dependencies. In J. F. Kolen and S. C. Kremer, editors, *A Field*
352 *Guide to Dynamical Recurrent Networks*. IEEE Press, 2001.
- 353 [21] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780,
354 1997.
- 355 [22] D. Horgan, J. Quan, D. Budden, G. Barth-Maron, M. Hessel, H. van Hasselt, and D. Silver.
356 Distributed prioritized experience replay. *ArXiv*, 2016. Sixth International Conference on
357 Learning Representations (ICLR).
- 358 [23] A. H. Klopff. Brain function and adaptive systems - a heterostatic theory. Technical Report
359 AFCRL-72-0164, Air Force Cambridge Research Laboratories, L. G. Hanscom Field, Bedford,
360 MA, 1972.
- 361 [24] J. Koutník, G. Cuccu, J. Schmidhuber, and F. Gomez. Evolving large-scale neural networks for
362 vision-based reinforcement learning. In *Proceedings of the 15th Annual Conference on Genetic*
363 *and Evolutionary Computation*, GECCO ’13, pages 1061–1068, 2013.
- 364 [25] W. Landecker, M. D. Thomure, L. M. A. Bettencourt, M. Mitchell, G. T. Kenyon, and S. P.
365 Brumby. Interpreting individual classifications of hierarchical networks. In *IEEE Symposium*
366 *on Computational Intelligence and Data Mining (CIDM)*, pages 32–38, 2013.
- 367 [26] S. Y. Lee, S. Choi, and S.-Y. Chung. Sample-efficient deep reinforcement learning via episodic
368 backward update. *ArXiv*, 2018.
- 369 [27] M. Lengyel and P. Dayan. Hippocampal contributions to control: The third way. In J. C. Platt,
370 D. Koller, Y. Singer, and S. T. Roweis, editors, *Advances in Neural Information Processing*
371 *Systems 20*, pages 889–896. Curran Associates, Inc., 2008.
- 372 [28] L. Li, T. J. Walsh, and M. L. Littman. Towards a unified theory of state abstraction for MDPs.
373 In *Ninth International Symposium on Artificial Intelligence and Mathematics (ISAIM)*, 2006.
- 374 [29] L. Lin. *Reinforcement Learning for Robots Using Neural Networks*. PhD thesis, Carnegie
375 Mellon University, Pittsburgh, 1993.
- 376 [30] J. Luoma, S. Ruutu, A. W. King, and H. Tikkanen. Time delays, competitive interdependence,
377 and firm performance. *Strategic Management Journal*, 38(3):506–525, 2017.
- 378 [31] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu.
379 Asynchronous methods for deep reinforcement learning. In M. F. Balcan and K. Q. Weinberger,
380 editors, *Proceedings of the 33rd International Conference on Machine Learning (ICML)*,
381 volume 48 of *Proceedings of Machine Learning Research*, pages 1928–1937. PMLR, 2016.
- 382 [32] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. A. Ried-
383 miller. Playing Atari with deep reinforcement learning. *ArXiv*, 2013.
- 384 [33] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
385 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou,
386 H. King, D. Kumaran, D. Wierstra, S. Legg, , and D. Hassabis. Human-level control through
387 deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

- 388 [34] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K.-R. Müller. Explaining nonlinear
389 classification decisions with deep Taylor decomposition. *Pattern Recognition*, 65:211 – 222,
390 2017.
- 391 [35] P. W. Munro. A dual back-propagation scheme for scalar reinforcement learning. In *Proceedings*
392 *of the Ninth Annual Conference of the Cognitive Science Society, Seattle, WA*, pages 165–176,
393 1987.
- 394 [36] A. Y. Ng, D. Harada, and S. J. Russell. Policy invariance under reward transformations: Theory
395 and application to reward shaping. In *Proceedings of the Sixteenth International Conference on*
396 *Machine Learning (ICML'99)*, pages 278–287, 1999.
- 397 [37] T. Pohlen, B. Piot, T. Hester, M. G. Azar, D. Horgan, D. Budden, G. Barth-Maron, H. van
398 Hasselt, J. Quan, M. Večerík, M. Hessel, R. Munos, and O. Pietquin. Observe and look further:
399 Achieving consistent performance on Atari. *ArXiv*, 2018.
- 400 [38] B. Poulin, R. Eisner, D. Szafron, P. Lu, R. Greiner, D. S. Wishart, A. Fyshe, B. Pearcy, C. Mac-
401 Donell, and J. Anvik. Visual explanation of evidence in additive classifiers. In *Proceedings*
402 *of the 18th Conference on Innovative Applications of Artificial Intelligence (IAAI)*, volume 2,
403 pages 1822–1829, 2006.
- 404 [39] M. L. Puterman. Markov decision processes. In *Stochastic Models*, volume 2 of *Handbooks in*
405 *Operations Research and Management Science*, chapter 8, pages 331–434. Elsevier, 1990.
- 406 [40] M. L. Puterman. *Markov Decision Processes*. John Wiley & Sons, Inc., 2005.
- 407 [41] H. Rahmandad, N. Repenning, and J. Sterman. Effects of feedback delay on learning. *System*
408 *Dynamics Review*, 25(4):309–338, 2009.
- 409 [42] B. Ravindran and A. G. Barto. Symmetries and model minimization in Markov decision
410 processes. Technical report, University of Massachusetts, Amherst, MA, USA, 2001.
- 411 [43] B. Ravindran and A. G. Barto. SMDP homomorphisms: An algebraic approach to abstraction in
412 semi-Markov decision processes. In *Proceedings of the 18th International Joint Conference on*
413 *Artificial Intelligence (IJCAI'03)*, pages 1011–1016, San Francisco, CA, USA, 2003. Morgan
414 Kaufmann Publishers Inc.
- 415 [44] A. J. Robinson. *Dynamic Error Propagation Networks*. PhD thesis, Trinity Hall and Cambridge
416 University Engineering Department, 1989.
- 417 [45] T. Robinson and F. Fallside. Dynamic reinforcement driven error propagation networks with
418 application to game playing. In *Proceedings of the 11th Conference of the Cognitive Science*
419 *Society, Ann Arbor*, pages 836–843, 1989.
- 420 [46] H. Sahni. Reinforcement learning never worked, and 'deep' only helped a bit. [himanshusahni.
421 github.io/2018/02/23/reinforcement-learning-never-worked.html](https://github.com/2018/02/23/reinforcement-learning-never-worked.html), 2018.
- 422 [47] T. Schaul, J. Quan, I. Antonoglou, and D. Silver. Prioritized experience replay. *ArXiv*, 2015.
- 423 [48] J. Schmidhuber. Making the world differentiable: On using fully recurrent self-supervised neural
424 networks for dynamic reinforcement learning and planning in non-stationary environments.
425 Technical Report FKI-126-90 (revised), Institut für Informatik, Technische Universität München,
426 1990. Experiments by Sepp Hochreiter.
- 427 [49] J. Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117,
428 2015.
- 429 [50] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
430 algorithms. *ArXiv*, 2018.
- 431 [51] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. van den Driessche, J. Schrit-
432 twieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham,
433 N. Kalchbrenner, I. Sutskever, T. P. Lillicrap, M. Leach, K. Kavukcuoglu, T. Graepel, and
434 D. Hassabis. Mastering the game of Go with deep neural networks and tree search. *Nature*,
435 529(7587):484–489, 2016.

- 436 [52] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre,
437 D. Kumaran, T. Graepel, T. P. Lillicrap, K. Simonyan, and D. Hassabis. Mastering Chess and
438 Shogi by self-play with a general reinforcement learning algorithm. *ArXiv*, 2017.
- 439 [53] S. P. Singh and R. S. Sutton. Reinforcement learning with replacing eligibility traces. *Machine*
440 *Learning*, 22:123–158, 1996.
- 441 [54] M. J. Sobel. The variance of discounted Markov decision processes. *Journal of Applied*
442 *Probability*, 19(4):794–802, 1982.
- 443 [55] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks. *ArXiv*, 2017.
- 444 [56] R. S. Sutton. Learning to predict by the methods of temporal differences. *Machine Learning*,
445 3:9–44, 1988.
- 446 [57] R. S. Sutton and A. G. Barto. Towards a modern theory of adaptive networks: expectation and
447 prediction. *Psychol. Rev.*, 88(2):135–170, 1981.
- 448 [58] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, Cambridge,
449 MA, 2 edition, 2017. Draft from November 2017.
- 450 [59] A. Tamar, D. DiCastro, and S. Mannor. Policy gradients with variance related risk criteria.
451 In J. Langford and J. Pineau, editors, *Proceedings of the 29th International Conference on*
452 *Machine Learning (ICML'12)*, 2012.
- 453 [60] H. van Hasselt. Double q -learning. In J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S.
454 Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems 23*, pages
455 2613–2621. Curran Associates, Inc., 2010.
- 456 [61] H. van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double q -learning.
457 In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, pages 2094–2100.
458 AAAI Press, 2016.
- 459 [62] Z. Wang, N. de Freitas, and M. Lanctot. Dueling network architectures for deep reinforcement
460 learning. *ArXiv*, 2015.
- 461 [63] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. de Freitas. Dueling network
462 architectures for deep reinforcement learning. In M. F. Balcan and K. Q. Weinberger, editors,
463 *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, volume 48 of
464 *Proceedings of Machine Learning Research*, pages 1995–2003. PMLR, 2016.
- 465 [64] E. Wiewiora. Potential-based shaping and q -value initialization are equivalent. *Journal of*
466 *Artificial Intelligence Research*, 19:205–208, 2003.
- 467 [65] E. Wiewiora, G. Cottrell, and C. Elkan. Principled methods for advising reinforcement learning
468 agents. In *Proceedings of the Twentieth International Conference on International Conference*
469 *on Machine Learning (ICML'03)*, pages 792–799, 2003.
- 470 [66] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
471 learning. *Machine Learning*, 8(3):229–256, 1992.
- 472 [67] Z. Xu, H. van Hasselt, and D. Silver. Meta-gradient reinforcement learning. *ArXiv*, 2018.
- 473 [68] J. Zhang, Z. L. Lin, J. Brandt, X. Shen, and S. Sclaroff. Top-down neural attention by excitation
474 backprop. In *Proceedings of the 14th European Conference on Computer Vision (ECCV)*, pages
475 543–559, 2016. part IV.