

1 **SPEEDY CATEGORICAL DISTRIBUTIONAL REINFORCEMENT**
2 **LEARNING AND COMPLEXITY ANALYSIS***

3 MARKUS BÖCK[†] AND CLEMENS HEITZINGER[‡]

4 **Abstract.** In distributional reinforcement learning, the entire distribution of the return instead
5 of just the expected return is modelled. The approach with categorical distributions as the approx-
6 imation method is well-known in Q-learning and convergence results have been established in the
7 tabular case. In this work, Speedy Q-learning is extended to categorical distributions, a finite-time
8 analysis is performed, and PAC bounds in terms of the Cramér distance are established. It is shown
9 that also in the distributional case the new update rule yields faster policy evaluation in comparison
10 to the standard Q-learning one and that the sample complexity is essentially the same as the one
11 of the value-based algorithmic counterpart. Without the need for more state-action-reward sam-
12 ples, one gains significantly more information about the return with categorical distributions. Even
13 though the results do not easily extend to the case of policy control, a slight modification to the
14 update rule yields promising numerical results.

15 **Key words.** Reinforcement learning, distributional reinforcement learning, Q-learning, PAC
16 bounds, complexity analysis.

17 **AMS subject classifications.** 68Q25 Analysis of algorithms and problem complexity, 68Q32
18 Computational learning theory, 68T05 Learning and adaptive systems in artificial intelligence, 68T42
19 Agent technology and artificial intelligence.

20 **1. Introduction.** Distributional reinforcement learning (DRL) is a subfield of
21 reinforcement learning where the entire return distribution is modelled directly, rather
22 than just the expected return. Bellemare et al. [1] introduced a particular distribu-
23 tional framework based on categorical distributions. Rowland et al. [8] established
24 convergence results in the tabular case for this approximation method. In 2011,
25 Ghavamzadeh et al. [5] introduced a new variant of Q-learning [11], called Speedy
26 Q-learning (SQL), which was subject to finite-time analysis and achieved impressive
27 experimental results. Finite-time analysis in terms of probably-approximately-correct
28 (PAC) bounds was also performed for Q-learning [3].

29 In this work, motivated by the results of Ghavamzadeh et al. [5] and Rowland et
30 al. [8], the SQL update rule is applied in the distributional framework and probably-
31 approximately-correct bounds for the resulting algorithm are established in terms of
32 the Cramér distance. It is shown that the sample complexity of this algorithm is
33 essentially the same as in the value-based case. Furthermore, the accelerated con-
34 vergence of SQL in terms of the expected return also translates to the distributional
35 case.

36 After presenting the theoretical background in Section 2, the Speedy Categorical
37 Policy Evaluation (SCPE) algorithm is introduced and the main theoretical results
38 are stated in Section 3. The corresponding proofs are given in Section 4. In Section 5,
39 the problems of using the SQL update rule in the control case are discussed. Lastly,
40 in Section 6, the theoretical results in the policy evaluation case are confirmed exper-
41 imentally, and it is shown that a slight modification to the update rule recovers the
42 improved convergence in the control case.

*Submitted to the editors 03.09.2020.

[†]Department of Mathematics and Geoinformation, TU Wien, Austria
(e1634838@student.tuwien.ac.at).

[‡]Department of Mathematics and Geoinformation, TU Wien, Austria
(clemens.heizinger@tuwien.ac.at).

2. Background.

2.1. Distributional Reinforcement Learning. The reinforcement learning objective is formalised by a Markov Decision Process (MDP), i.e., a tuple $\langle \mathcal{X}, \mathcal{A}, r, p \rangle$, where \mathcal{X} is a set of states and \mathcal{A} is a set of actions. In this work, we only consider finite MDPs, i.e., $|\mathcal{X}| < \infty$ and $|\mathcal{A}| < \infty$. Trajectories (X_t, A_t, R_t) are obtained through the selection of actions at given states, where the probabilities of state transitions are defined by the deterministic function p and R_t are defined by the kernel r , where $r(\cdot|x, a, x')$ represents the immediate reward when transitioning from state x with action a to state x' . The Markov property requires that X_{t+1} and R_t only depend on the previous state and action (x, a) , i.e.,

$$\mathbb{P}(R_t = s, X_{t+1} = x' | X_t = x, A_t = a, X_{t-1} = x_{t-1}, \dots) = r(s|x, a, x')p(x'|x, a).$$

The standard approach to reinforcement learning is to model the expected return, usually by a state-action value function Q . As the name suggests, the core of distributional reinforcement learning is to model the entire distribution of the return directly.

For $(x, a) \in \mathcal{X} \times \mathcal{A}$, the *return* $Z^\pi(x, a)$ is the sum of discounted rewards along a trajectory following the policy π starting in state x and taking action a , i.e.,

$$Z^\pi(x, a) := \sum_{t=0}^{\infty} \gamma^t R_t,$$

$$\text{where } X_0 = x, \quad A_0 = a, \quad X_{t+1} \sim p(\cdot|X_t, A_t),$$

$$A_{t+1} \sim \pi(\cdot|X_{t+1}), \quad R_t \sim r(\cdot|X_t, A_t, X_{t+1}).$$

The function Z^π mapping state-action pairs to random variables is called *return distribution function*.

The usual state-action value function Q^π can be related to the return distribution function by observing that

$$Q^\pi(x, a) = \mathbb{E}[Z^\pi(x, a)].$$

Furthermore, the Bellman equation [2] can be extended to the distributional case as

$$(2.1) \quad Z^\pi(x, a) \stackrel{D}{=} R + \gamma Z^\pi(X', A'),$$

where $X' \sim p(\cdot|x, a)$, $A' \sim \pi(\cdot|X')$ and $R \sim r(\cdot|x, a, X')$. Here the equal sign indicates that the random variable on the left hand side and the one on the right hand side are identically distributed.

Let $\eta_\pi^{(x,a)}$ denote the underlying probability distribution of the random variable $Z^\pi(x, a)$, giving us a second representation

$$Z^\pi(x, a) \sim \eta_\pi^{(x,a)}$$

of return distribution functions in terms of probability measures.

Let $\eta: \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{P}(\mathbb{R})$ be an arbitrary mapping to probability measures on \mathbb{R} . For a given policy π , the distributional Bellman operator \mathcal{T}^π can be written in terms of cumulative distribution functions as

$$(2.2) \quad F_{\mathcal{T}^\pi \eta^{(x,a)}}(z) = \mathbb{E} \left[F_{\eta^{(X', A')}} \left(\frac{z - R}{\gamma} \right) \right]$$

due to (2.1). If η corresponds to the return distributions of π , i.e., $\eta = \eta_\pi$, it follows directly from the Bellman equation that $\mathcal{T}^\pi \eta_\pi = \eta_\pi$.

84 **2.2. Categorical DRL.** A challenge of distributional reinforcement learning is
 85 to find appropriate methods to approximate the return distributions. Bellemare et
 86 al. [1] proposed to use N fixed atoms z_1, \dots, z_N or grid points and defined the set of
 87 categorical distributions as

$$88 \quad \mathcal{P}_z := \left\{ \sum_{i=1}^N p_i \delta_{z_i} \mid p_i \geq 0 \wedge \sum_{i=1}^N p_i = 1 \right\}.$$

89 As this set is not closed under the Bellman update, we have to project distributions
 90 back onto this set. Thus the categorical projection operator $\Pi_{\mathcal{C}}$ was introduced, which
 91 is explained in more detail later.

92 Rowland et al. [8] connected this operator to the Cramér distance, which is
 93 defined between two return distribution functions as

$$94 \quad \bar{\ell}_2(\eta, \xi) := \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \ell_2(\eta^{(x,a)}, \xi^{(x,a)})$$

$$95 \quad = \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \left(\int_{\mathbb{R}} |F_{\eta^{(x,a)}}(z) - F_{\xi^{(x,a)}}(z)|^2 dz \right)^{1/2}.$$

96

97 The key observation was that

$$98 \quad \Pi_{\mathcal{C}} \mathcal{T}^{\pi} : \mathcal{P}_z^{\mathcal{X} \times \mathcal{A}} \rightarrow \mathcal{P}_z^{\mathcal{X} \times \mathcal{A}}$$

99 is a $\sqrt{\gamma}$ -contraction in $\bar{\ell}_2$. With this fact, the convergence of categorical policy eval-
 100 uation to $\eta_{\mathcal{C}}$, the unique fixed point of $\Pi_{\mathcal{C}} \mathcal{T}^{\pi}$, was proven. Note that since we only
 101 approximate distributions, we have $\eta_{\mathcal{C}} \neq \eta_{\pi}$ in general. However, if the return dis-
 102 tributions $\eta_{\pi}^{(x,a)}$ are supported on $[z_1, z_N]$ for all $(x, a) \in \mathcal{X} \times \mathcal{A}$, then increasing the
 103 number of atoms yields a better precision, i.e.,

$$104 \quad \bar{\ell}_2^2(\eta_{\mathcal{C}}, \eta_{\pi}) \leq \frac{1}{1-\gamma} \max_{1 \leq i < N} (z_{i+1} - z_i).$$

105 For policy evaluation, Rowland et al. [8] considered the update rule given by the
 106 weighted sum

$$107 \quad (2.3) \quad \eta_{k+1}^{(x,a)} := (1 - \alpha_k(x, a)) \eta_k^{(x,a)} + \alpha_k(x, a) \Pi_{\mathcal{C}} \mathcal{T}_k^{\pi} \eta_k^{(x,a)}.$$

108 Here \mathcal{T}_k^{π} is the stochastic Bellman operator at time k , which depends on samples
 109 $x'_k \sim p(\cdot | x, a)$ and $a'_k \sim \pi(\cdot | x'_k)$ as well as the reward sample $r_k \sim r(\cdot | x, a, x'_k)$ for each
 110 $(x, a) \in \mathcal{X} \times \mathcal{A}$. In terms of cumulative distribution functions, the operator can be
 111 written as

$$112 \quad (2.4) \quad F_{\mathcal{T}_k^{\pi} \eta^{(x,a)}}(z) = F_{\eta^{(x'_k, a'_k)}} \left(\frac{z - r_k}{\gamma} \right),$$

113 which is a random variable for all $z \in \mathbb{R}$ due to (2.1), and we have

$$114 \quad F_{\mathcal{T}^{\pi} \eta^{(x,a)}}(z) = \mathbb{E} \left[F_{\mathcal{T}_k^{\pi} \eta^{(x,a)}}(z) \right].$$

115 In the update rule, $\alpha_k(x, a)$ are stepsizes. If $\eta_{\pi}^{(x,a)}$ is supported on $[z_1, z_N]$ and
 116 the Robins-Monro conditions $\sum_{k=0}^{\infty} \alpha_k(x, a) = \infty$ and $\sum_{k=0}^{\infty} \alpha_k(x, a)^2 < \infty$ hold for

117 all $(x, a) \in \mathcal{X} \times \mathcal{A}$, then the Cramér distance to the fixed point of $\Pi_C \mathcal{T}^\pi$ converges to
 118 zero almost surely, i.e., $\bar{\ell}_2(\eta_k, \eta_C) \rightarrow 0$ [8, Theorem 1].

119 In practice, the state and action samples usually come in episodes and at time k ,
 120 η_k is only updated at the current state-action pair (x_k, a_k) in the trajectory, and
 121 we have $\alpha_k(x, a) = 0$ for all $(x, a) \neq (x_k, a_k)$. This method presents the categorical
 122 distributional analog to the temporal difference (TD) algorithm [9, 10].

123 For policy control, we change \mathcal{T}_k^π in (2.3) to the stochastic optimality operator
 124 \mathcal{T}_k given by

$$125 \quad (2.5) \quad F_{\mathcal{T}_k \eta(x,a)}(z) = F_{\eta(x'_k, a_k^*)} \left(\frac{z - r_k}{\gamma} \right), \quad a_k^* = \operatorname{argmax}_{a \in \mathcal{A}} \mathbb{E}_{Z \sim \eta_k(x'_k, a)} [Z],$$

126 and obtain the categorical distributional equivalent to Q-learning [11, 10]. In this
 127 control case, convergence was only established with the additional assumption that a
 128 unique optimal policy exists [8, Theorem 2].

129 **2.3. Speedy Q-learning and Sample Complexities.** For a sample x, a, r, x' ,
 130 where $x' \sim p(\cdot | x, a)$, the Q-learning [11] update rule reads

$$131 \quad (2.6) \quad Q_{k+1}(x, a) = (1 - \alpha_k(x, a))Q_k(x, a) + \alpha_k(x, a)(r + \gamma \max_{a \in \mathcal{A}} Q_k(x', a)).$$

132 Let Q^* denote the unique optimal value function. Further, assume that the
 133 rewards are bounded by R_{\max} . For $\gamma < 1$ and $\beta := \frac{1}{1-\gamma}$, let $V_{\max} := \beta R_{\max}$ be the
 134 maximal attainable return.

135 If the updates with (2.6) are performed *synchronously*, that is at each time step k
 136 all state-action pairs $(x, a) \in \mathcal{X} \times \mathcal{A}$ are updated, and we have polynomial learning
 137 rates

$$138 \quad \alpha_k = \frac{1}{(k+1)^\omega}, \quad \frac{1}{2} < \omega < 1,$$

139 then for a finite state-action space $n = |\mathcal{X} \times \mathcal{A}|$, and for $\gamma < 1$, the following finite-time
 140 behaviour is known [3]: with probability at least $1 - \delta$, the inequality

$$141 \quad (2.7) \quad \|Q^* - Q_T\|_\infty \leq \epsilon$$

142 holds for

$$143 \quad T \geq C \left(\left(\frac{\beta^4 R_{\max}^2 \log \frac{n\beta^2 R_{\max}}{\delta\epsilon}}{\epsilon^2} \right)^{\frac{1}{\omega}} + \left(\beta \log \frac{\beta R_{\max}}{\epsilon} \right)^{\frac{1}{1-\omega}} \right)$$

144 and for some constant $C > 0$.

145 Following the reasoning of Even-Dar et al. [3] and Ghavamzadeh et al. [5], if γ
 146 is close to 1, β becomes the dominant term and the bound is optimized for $\omega = 4/5$,
 147 yielding a complexity of

$$148 \quad \mathcal{O} \left(\left(\frac{\beta^4 R_{\max}^2 \log \frac{n\beta^2 R_{\max}}{\delta\epsilon}}{\epsilon^2} \right)^{5/4} \right) = \tilde{\mathcal{O}}(\beta^5 / \epsilon^{2.5}),$$

149 since $g = \tilde{\mathcal{O}}(f) \iff g \leq C_1 f \log^{C_2}(f)$ for some constants $C_1, C_2 > 0$.

150 The bound in probability and the derived sample complexity also hold for the
 151 evaluation of the value function Q^π for an arbitrary policy π .

152 Ghavamzadeh et al. [5] introduced a faster variant of Q-learning and coined it
 153 Speedy Q-learning (SQL). They defined the update rule

(2.8)

$$154 Q_{k+1}(x, a) := (1 - \alpha_k)Q_k(x, a) + \alpha_k \mathcal{T}_k Q_{k-1}(x, a) + (1 - \alpha_k)(\mathcal{T}_k Q_k(x, a) - \mathcal{T}_k Q_{k-1}(x, a))$$

155 based on two previous time steps instead of just one, where

$$156 \mathcal{T}_k Q(x, a) = r + \gamma \max_{a \in \mathcal{A}} Q(x', a)$$

157 and the learning rate is $\alpha_k = \frac{1}{k+1}$. The key difference to Q-learning is that SQL uses
 158 a more aggressive learning rate for the third term. Changing it to $\alpha_k(\mathcal{T}_k Q_k(x, a) -$
 159 $\mathcal{T}_k Q_{k-1}(x, a))$ would be equivalent to Q-learning. The difference seems small, however
 160 it yields faster convergence, i.e., it can be shown that the inequality

$$161 (2.9) \quad \|Q^* - Q_T\|_\infty \leq 2V_{\max}\beta \left(\frac{\gamma}{T} + \sqrt{\frac{2 \log \frac{2n}{\delta}}{T}} \right)$$

162 holds with probability $1 - \delta$ and thus for

$$163 T := \frac{11.66\beta^2 V_{\max}^2 \log \frac{2n}{\delta}}{\epsilon^2}.$$

164 we have

$$165 \|Q^* - Q_T\|_\infty \leq \epsilon.$$

166 Again, viewing β as the dominant term, we have a convergence rate of $\tilde{O}(\beta^4/\epsilon^2)$.

167 **3. Speedy Categorical Policy Evaluation (SCPE).** In the following, the
 168 update rule of Speedy Q-learning is extended to categorical distributions in the policy
 169 evaluation case. We chose to extend SQL to distributions, rather than standard Q-
 170 learning, because it yields faster convergence. However, it is worth mentioning that
 171 the main idea of the proof is also applicable if one uses (2.3) and (2.6).

172 In order to translate SQL to categorical distributions, we combine (2.8) and (2.3)
 173 for the evaluation of a fixed policy π into the update formula

$$174 (3.1) \quad \eta_{k+1}^{(x,a)} = \eta_k^{(x,a)} + \alpha_k (\Pi_C \mathcal{T}_k^\pi \eta_{k-1}^{(x,a)} - \eta_k^{(x,a)}) + (1 - \alpha_k) (\Pi_C \mathcal{T}_k^\pi \eta_k^{(x,a)} - \Pi_C \mathcal{T}_k^\pi \eta_{k-1}^{(x,a)}),$$

175 where we start with two initial return distribution functions $\eta_0 = \eta_{-1} \in \mathcal{P}_z$. We again
 176 use the learning rate $\alpha_k := \frac{1}{k+1}$.

177 It is straightforward to see that (3.1) can be rewritten as the convex combination

$$178 \eta_{k+1}^{(x,a)} = \frac{k}{k+1} \eta_k^{(x,a)} + \frac{1}{k+1} \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)},$$

179 where we define the sample update as

$$180 \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} := k \Pi_C \mathcal{T}_k^\pi \eta_k^{(x,a)} - (k-1) \Pi_C \mathcal{T}_k^\pi \eta_{k-1}^{(x,a)}.$$

181 Note that it is ad hoc not clear whether $\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)}$ is a probability measure. In
 182 general, it is a finite signed measure, and thus we also do not know if the recursively
 183 defined $\eta_k^{(x,a)}$ are indeed probability measures. The consideration of this problem
 184 makes up a substantial part of the analysis in Section 4.

185 In the following analysis, we only consider the synchronous version of policy eval-
 186 uation, which is shown as pseudo-code in Algorithm 3.1. Like the finite-time analysis
 187 of SQL and Q-learning, it can also be extended to the asynchronous case, where we
 188 consider a policy with finite covering time.

Algorithm 3.1 Synchronous Speedy Categorical Policy Evaluation

```

1: Input: discount factor  $\gamma$ , policy  $\pi$ , number of iterations  $T$ , initial guess  $\eta_0$ 
2:  $\eta_{-1} \leftarrow \eta_0$ 
3: for  $k \in 0, \dots, T - 1$  do
4:    $\alpha_k \leftarrow \frac{1}{k+1}$ 
5:   for  $(x, a) \in \mathcal{X} \times \mathcal{A}$  do
6:     Sample  $x'_k \sim p(\cdot|x, a)$ ,  $a'_k \sim \pi(\cdot|x'_k)$ ,  $r_k \sim r(\cdot|x, a, x'_k)$ 
7:      $\mathcal{T}_k^\pi \eta_k^{(x,a)} \leftarrow \sum_{i=1}^N p_{k,i}^{(x'_k, a'_k)} \delta_{r_k + \gamma z_i}$  # Bellman update
8:      $\mathcal{T}_k^\pi \eta_{k-1}^{(x,a)} \leftarrow \sum_{i=1}^N p_{k-1,i}^{(x'_k, a'_k)} \delta_{r_k + \gamma z_i}$  # Bellman update
9:     # Project onto support  $z_1, \dots, z_N$  and calculate difference
10:     $\mathcal{D}_k^{(x,a)} \leftarrow k \Pi_{\mathcal{C}} \mathcal{T}_k^\pi \eta_k^{(x,a)} - (k-1) \Pi_{\mathcal{C}} \mathcal{T}_k^\pi \eta_{k-1}^{(x,a)}$ 
11:    # Update  $\eta$ 
12:     $\eta_{k+1}^{(x,a)} \leftarrow (1 - \alpha_k) \eta_k^{(x,a)} + \alpha_k \mathcal{D}_k^{(x,a)}$ 
13:   end for
14: end for

```

189 In order to formulate the main result below, we collect the following assumptions.

190 ASSUMPTION 1. *The state-action space is finite with $n := |\mathcal{X} \times \mathcal{A}|$ elements.*
 191 *The categorical distribution $\eta_{\mathcal{C}}$ is the unique fixed point of $\Pi_{\mathcal{C}} \mathcal{T}^\pi$. The rewards are*
 192 *bounded by $R_{\max} > 0$. The discount factor γ is smaller than 1, and let $\bar{\beta} := \frac{1}{1-\sqrt{\gamma}}$. Let*
 193 *$V_{\max} := \frac{1}{1-\gamma} R_{\max}$ be the maximal attainable return. For the N fixed atoms we assume*
 194 *$z_1 = -V_{\max}$ and $z_N = V_{\max}$. Lastly, the two initial return distribution functions are*
 195 *equal, i.e., $\eta_{-1} = \eta_0$, and the η_k are obtained by update rule (3.1).*

196 The main result is the following.

197 THEOREM 3.1. *Under Assumption 1, the inequality*

$$198 \quad \bar{\ell}_2(\eta_{\mathcal{C}}, \eta_T) \leq \sqrt{2V_{\max} \bar{\beta}} \left(\frac{\sqrt{\gamma}}{T} + \sqrt{\frac{2 \log \frac{2nN}{\delta}}{T}} \right)$$

199 *holds with probability at least $1 - \delta$.*

200 We give two corollaries.

201 COROLLARY 3.2. *Under Assumption 1, for any $0 < \epsilon \leq \sqrt{V_{\max}}$, the inequality*
 202 *$\|\eta_{\mathcal{C}} - \eta_T\|_{\bar{\ell}_2} \leq \epsilon$ holds with probability at least $1 - \delta$ after*

$$203 \quad T := \frac{6.53 \bar{\beta}^2 V_{\max} \log \frac{2nN}{\delta}}{\epsilon^2}$$

204 *steps of SCPE.*

205 COROLLARY 3.3. *Under Assumption 1, η_T converges to $\eta_{\mathcal{C}}$ almost surely in $\bar{\ell}_2$.*

206 The proofs are deferred to Section 4.

207 Corollary 3.2 leads to following complexity analysis. For each time step k , we
 208 sweep over the entire state-action space. Therefore, after T iterations, $3nT$ samples
 209 are available in total (reward, next state, and next action in each time step). For
 210 γ close to 1, we have $\tilde{\beta} \approx 2\beta$. Recall that $V_{\max} = \beta R_{\max}$. Therefore, the sample
 211 complexity of SCPE is

$$212 \quad \tilde{O}(n\beta^3/\epsilon^2)$$

213 (omitting the logarithmic factor). The number N of atoms only contributes to the
 214 logarithmic factor. Thus, increasing the accuracy of the distribution approximation
 215 causes only a small penalty.

216 Further, SCPE has *essentially the same* sample complexity as value-based SQL,
 217 which is

$$218 \quad \tilde{O}(n\beta^4/\epsilon^2).$$

219 The difference in the power of β stems from the fact that a different metric was used.
 220 To see how the difference in expected values and the Cramér distance relate, consider
 221 two measures μ, ν supported on $[z_1, z_N]$. Then,

$$\begin{aligned} 222 \quad & |\mathbb{E}_{Z_\mu \sim \mu} [Z_\mu] - \mathbb{E}_{Z_\nu \sim \nu} [Z_\nu]| = \\ 223 \quad & = \left| \int_0^\infty (1 - F_\mu(z)) - (1 - F_\nu(z)) dz - \int_{-\infty}^0 F_\mu(z) - F_\nu(z) dz \right| \\ 224 \quad & \leq \int_{\mathbb{R}} |F_\mu(z) - F_\nu(z)| dz \\ 225 \quad & \leq \|F_\mu - F_\nu\|_2 \|\mathbf{1}_{[z_1, z_N]}\|_2 \\ 226 \quad & = (z_N - z_1)^{1/2} \|F_\mu - F_\nu\|_2 = \sqrt{2V_{\max}} \ell_2(\mu, \nu) \end{aligned}$$

228 This inequality precisely captures the relationship of inequality (2.9) and Theorem
 229 3.1. As $V_{\max} = \beta R_{\max}$ this also explains the difference in the power of β in the sample
 230 complexities.

231 It is quite an interesting result that the sample complexity remains the same
 232 when switching to distributions. One *does not* need more samples when modelling
 233 the entire distribution. Further, the sample complexity is independent of the number
 234 of atoms - the precision with which the return distributions are modelled. However,
 235 the computational complexity $\tilde{O}(nN\beta^3/\epsilon^2)$ is higher, of course, and a table with nN
 236 elements is needed to store the return distributions.

237 **4. Analysis.** The analysis follows the outline of Ghavamzadeh et al. [5]. Since
 238 in DRL the return distributions depend on state, action, and reward samples, it is
 239 imperative to extend the notion of random variables to random distributions. We
 240 define signed random measures according to Kallenberg [6].

241 DEFINITION 4.1. *Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and define*

$$242 \quad M := \{ \nu \text{ signed measure on } (\mathbb{R}, \mathcal{B}) \mid |\nu(B)| < \infty \text{ for all bounded } B \in \mathcal{B} \},$$

243 where \mathcal{B} is the Borel- σ -field on \mathbb{R} . M is equipped with the σ -field \mathcal{M} , which is the
 244 smallest σ -field such that $\nu \mapsto \nu(B)$ is measurable for all $B \in \mathcal{B}$.

245 Measurable functions $X: (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (M, \mathcal{M})$, $\omega \mapsto X_\omega$, are called signed random
 246 measures.

247 The expected measure $\mathbb{E}[X] \in M$ is given by

$$248 \quad \mathbb{E}[X](A) := \mathbb{E}[X(A)], \quad \text{where } X(A): \Omega \rightarrow \mathbb{R}, \quad \omega \mapsto X_\omega(A).$$

249 Further, $F_X(z) := (\omega \mapsto F_{X_\omega}(z))$ is a random variable for all $z \in \mathbb{R}$, and we have

$$250 \quad (4.1) \quad F_{\mathbb{E}[X]}(z) = \mathbb{E}[X]((-\infty, z]) = \mathbb{E}[X(-\infty, z]] = \mathbb{E}[F_X(z)].$$

251 The set of all signed random measures on $E \subseteq M$ is denoted by

$$252 \quad \mathcal{P}(E) := \{f : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (E, \mathcal{M}|_E) \text{ measurable}\}.$$

253 **4.1. Step 1: Stability.** As mentioned, we do not know whether $\eta_k^{(x,a)}$ are indeed
254 probability measures. For that reason, we first define a vector space of finite signed
255 measures, which allows us to freely perform addition and scalar multiplication.

256 DEFINITION 4.2. Let \mathcal{L} be the set of finite signed Borel measures

$$257 \quad \mathcal{L} = \{\nu \text{ signed measure} \mid \exists F_\nu : \mathbb{R} \rightarrow \mathbb{R} \text{ right continuous,} \\ 258 \quad \nu((a, b]) = F_\nu(b) - F_\nu(a), \quad |\nu(\mathbb{R})| < \infty, \quad \lim_{z \rightarrow -\infty} F_\nu(z) = 0, \quad \lim_{z \rightarrow \infty} F_\nu(z) < \infty\}$$

260 \mathcal{L} becomes a real vector space by defining

$$261 \quad (4.2) \quad (a\mu + b\nu)(A) := a\mu(A) + b\nu(A), \quad \mu, \nu \in \mathcal{L}, \quad a, b \in \mathbb{R}, \quad A \text{ a measurable set.}$$

262 Equation (4.2) immediately implies that

$$263 \quad (4.3) \quad F_{a\mu + b\nu} = aF_\mu + bF_\nu.$$

264 The categorical distributions are also extended to a subspace of signed measures,

$$265 \quad \mathcal{P}_z \subseteq \mathcal{L}_z := \left\{ \sum_{i=1}^N c_i \delta_{z_i} \mid c_i \in \mathbb{R} \right\} \subseteq \mathcal{L}.$$

266 The categorical projection operator Π_C can be easily applied to elements of \mathcal{L} by
267 defining

$$268 \quad (4.4) \quad \Pi_C : \mathcal{L} \rightarrow \mathcal{L}_z, \quad F_{\Pi_C \nu}(z_i) = \frac{1}{z_{i+1} - z_i} \int_{z_i}^{z_{i+1}} F_\nu(z) dz, \quad F_{\Pi_C \nu}(z_N) = \lim_{z \rightarrow \infty} F_\nu(z).$$

269 From (4.4) and (4.3), it is not difficult to see that $\Pi_C : \mathcal{L} \rightarrow \mathcal{L}_z$ is a linear
270 projection. Furthermore, from characterisation (2.4) and (4.3) it follows that also
271 $\mathcal{T}_k^\pi : \mathcal{L}^{\mathcal{X} \times \mathcal{A}} \rightarrow \mathcal{L}^{\mathcal{X} \times \mathcal{A}}$ is a linear mapping.

272 Recall that $\mathcal{P}(\mathcal{P}_z)$ in the next lemma is the set of random measures with values
273 in \mathcal{P}_z .

274 LEMMA 4.3. For all $k \geq 0$, it holds that $\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} \in \mathcal{P}(\mathcal{P}_z)$ and $\eta_k^{(x,a)} \in$
275 $\mathcal{P}(\mathcal{P}_z)$.

276 *Proof.* This result is proved by induction. Since we only extended $\Pi_C \mathcal{T}_k^\pi$ to signed
277 measures, it is still true that when passed a (random) probability measure $\Pi_C \mathcal{T}_k^\pi$
278 outputs a random probability measure.

279 Recall that $\mathcal{D}_k[\eta_k, \eta_{k-1}] = k\Pi_C \mathcal{T}_k^\pi \eta_k - (k-1)\Pi_C \mathcal{T}_k^\pi \eta_{k-1}$. As the initial return
280 distributions are identical, we have

$$281 \quad \mathcal{D}_0[\eta_0, \eta_{-1}]^{(x,a)} = \Pi_C \mathcal{T}_0^\pi \eta_{-1}^{(x,a)} = \Pi_C \mathcal{T}_0^\pi \eta_0^{(x,a)}.$$

282 $\mathcal{D}_0[\eta_k, \eta_{k-1}]^{(x,a)}$ is a random probability measure and an element of $\mathcal{P}(\mathcal{P}_z)$, since
 283 $\eta_0^{(x,a)} \in \mathcal{P}_z$. Of course, $\eta_0^{(x,a)} \in \mathcal{P}(\mathcal{P}_z)$ also (interpreted as a random measure which
 284 takes $\eta_0^{(x,a)}$ with probability 1).

285 Assume that $\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)}$ and $\eta_k^{(x,a)}$ are random probability measures. To
 286 show the induction step, we can relate $\mathcal{D}_{k+1}[\eta_{k+1}, \eta_k]$ to $\mathcal{D}_k[\eta_k, \eta_{k-1}]$ by observing
 287 that

$$\begin{aligned}
 & \mathcal{D}_{k+1}[\eta_{k+1}, \eta_k]^{(x,a)} \\
 &= (k+1)\Pi_C \mathcal{T}_{k+1}^\pi \eta_{k+1}^{(x,a)} - k\Pi_C \mathcal{T}_{k+1}^\pi \eta_k^{(x,a)} \\
 &= (k+1)\Pi_C \mathcal{T}_{k+1}^\pi \left(\frac{k}{k+1} \eta_k + \frac{1}{k+1} \mathcal{D}_k[\eta_k, \eta_{k-1}] \right)^{(x,a)} - k\Pi_C \mathcal{T}_{k+1}^\pi \eta_k^{(x,a)} \\
 &= k\Pi_C \mathcal{T}_{k+1}^\pi \eta_k^{(x,a)} + \Pi_C \mathcal{T}_{k+1}^\pi \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} - k\Pi_C \mathcal{T}_{k+1}^\pi \eta_k^{(x,a)} \\
 &= \Pi_C \mathcal{T}_{k+1}^\pi \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)},
 \end{aligned}$$

294 where we used the fact that $\Pi_C \mathcal{T}_k^\pi$ is linear. Thus, $\mathcal{D}_{k+1}[\eta_{k+1}, \eta_k]^{(x,a)} \in \mathcal{P}(\mathcal{P}_z)$ also.

295 Since

$$\eta_{k+1} = \frac{k}{k+1} \eta_k + \frac{1}{k+1} \mathcal{D}_k[\eta_k, \eta_{k-1}]$$

297 and \mathcal{P}_z is a convex set, we have $\eta_{k+1}^{(x,a)} \in \mathcal{P}(\mathcal{P}_z)$. \square

298 **4.2. Step 2: Error Martingale.** The history of the algorithm at time k can
 299 be captured in the form of the filtration

$$\mathcal{F}_k := \sigma\text{-field generated by } r_1, x'_1, a'_1, \dots, r_k, x'_k, a'_k, \quad (x, a) \in \mathcal{X} \times \mathcal{A}.$$

300 The expected update is given by

$$\mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)} := \mathbb{E} \left[\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} \mid \mathcal{F}_{k-1} \right] \stackrel{(2.4)}{=} k\Pi_C \mathcal{T}^\pi \eta_k^{(x,a)} - (k-1)\Pi_C \mathcal{T}^\pi \eta_{k-1}^{(x,a)}.$$

303 The error $\epsilon_k^{(x,a)}$ and the cumulative error to the sample update $E_k^{(x,a)}$ are given by

$$\begin{aligned}
 \epsilon_k^{(x,a)} &:= \mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)} - \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)}, \\
 E_k^{(x,a)} &:= \sum_{j=0}^k \epsilon_j^{(x,a)}.
 \end{aligned}$$

307 Again, we can rewrite the update rule in terms of the expected update and the error
 308 as

$$(4.5) \quad \eta_{k+1}^{(x,a)} = \frac{k}{k+1} \eta_k^{(x,a)} + \frac{1}{k+1} (\mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)} - \epsilon_k^{(x,a)}).$$

310 It is not immediately clear how one can turn the errors into a martingale. The
 311 following Lemma shows that we have to look at the cumulative distribution function
 312 at each atom. Lemma 4.3 and Lemma 4.4 are the core results that allow us to extend
 313 the analysis of Speedy Q-learning [5] to categorical distributions. One can extend the
 314 result (2.7) from Even-Dar et al. [3] in a similar fashion.

315 LEMMA 4.4. *The inclusions $\epsilon_k^{(x,a)} \in \mathcal{P}(\mathcal{L}_z)$ and $E_k^{(x,a)} \in \mathcal{P}(\mathcal{L}_z)$ hold for all*
 316 *$k \geq 0$. For each atom z_i , it holds that the cumulative distribution functions of the*
 317 *error ϵ_k evaluated at z_i form a uniformly bounded martingale difference sequence, i.e.,*

$$318 \quad (4.6) \quad \forall k \geq 0: \quad \mathbb{E} \left[F_{\epsilon_k^{(x,a)}}(z_i) \mid \mathcal{F}_{k-1} \right] = 0 \wedge \left| F_{\epsilon_k^{(x,a)}}(z_i) \right| \leq 1.$$

319 *Proof.* By Lemma 4.3, $\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} \in \mathcal{P}(\mathcal{P}_z)$ holds. It follows from (4.1)
 320 that the expected measure $\mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)} \in \mathcal{P}_z$. This makes $\epsilon_k^{(x,a)}$ the difference of
 321 a random probability measure in $\mathcal{P}(\mathcal{P}_z)$ and a probability measure in \mathcal{P}_z . Therefore
 322 it is an element of $\mathcal{P}(\mathcal{L}_z)$. Further, $E_k^{(x,a)}$ is the sum of elements of $\mathcal{P}(\mathcal{L}_z)$ and thus
 323 also in $\mathcal{P}(\mathcal{L}_z)$.

324 By definition,

$$325 \quad \mathbb{E} \left[\epsilon_k^{(x,a)} \mid \mathcal{F}_{k-1} \right] = \mathbb{E} \left[\mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)} - \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} \mid \mathcal{F}_{k-1} \right]$$

$$326 \quad = \mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)} - \mathbb{E} \left[\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} \mid \mathcal{F}_{k-1} \right] = 0 \in \mathcal{L}_z,$$

328 and therefore

$$329 \quad \mathbb{E} \left[F_{\epsilon_k^{(x,a)}}(z_i) \mid \mathcal{F}_{k-1} \right] = F_{\mathbb{E}[\epsilon_k^{(x,a)} \mid \mathcal{F}_{k-1}]}(z_i) = 0 \in \mathbb{R}.$$

330 Furthermore, we have that

$$332 \quad F_{\epsilon_k^{(x,a)}}(z_i) = F_{\mathcal{D}[\eta_k, \eta_{k-1}]^{(x,a)}}(z_i) - F_{\mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)}}(z_i)$$

333 is the difference of a real value in $[0, 1]$ and a random variable with values in $[0, 1]$.

334 This makes it a random variable which is bounded by 1. \square

335 **4.3. Step 3: Upper bound.** The following lemma shows that $\eta_k \approx \Pi_C \mathcal{T}^\pi \eta_{k-1}$.

336 LEMMA 4.5. *For all $k \geq 1$, the equality*

$$337 \quad \eta_k = \frac{1}{k} (\Pi_C \mathcal{T}^\pi \eta_0 + (k-1) \Pi_C \mathcal{T}^\pi \eta_{k-1} - E_{k-1})$$

338 holds.

339 *Proof.* The equation is proved by induction. The result holds for $k = 1$, since

$$340 \quad \eta_1 = \mathcal{D}[\eta_0, \eta_{-1}] - \epsilon_0 = \Pi_C \mathcal{T}^\pi \eta_{-1} - \epsilon_0 = \Pi_C \mathcal{T}^\pi \eta_0 - E_0.$$

341 Assume that the equation holds for $k \geq 1$. The definitions of $\mathcal{D}[\eta_k, \eta_{k-1}]$ and E_k
 342 imply

$$343 \quad \eta_{k+1}$$

$$344 \quad = \frac{k}{k+1} \eta_k + \frac{1}{k+1} (\mathcal{D}[\eta_k, \eta_{k-1}] - \epsilon_k)$$

$$345 \quad = \frac{k}{k+1} \eta_k + \frac{1}{k+1} (k \Pi_C \mathcal{T}^\pi \eta_k - (k-1) \Pi_C \mathcal{T}^\pi \eta_{k-1} - \epsilon_k)$$

$$346 \quad = \frac{k}{k+1} \left(\frac{1}{k} (\Pi_C \mathcal{T}^\pi \eta_0 + (k-1) \Pi_C \mathcal{T}^\pi \eta_{k-1} - E_{k-1}) \right)$$

$$347 \quad \quad \quad + \frac{1}{k+1} (k \Pi_C \mathcal{T}^\pi \eta_k - (k-1) \Pi_C \mathcal{T}^\pi \eta_{k-1} - \epsilon_k)$$

$$348 \quad = \frac{1}{k+1} (\Pi_C \mathcal{T}^\pi \eta_0 + k \Pi_C \mathcal{T}^\pi \eta_k - E_{k-1} - \epsilon_k) = \frac{1}{k+1} (\Pi_C \mathcal{T}^\pi \eta_0 + k \Pi_C \mathcal{T}^\pi \eta_k - E_k),$$

350 which concludes the proof. \square

351 As \mathcal{L}_z is a vector space, it is more convenient to work with norms instead of
 352 metrics. For that matter, we define

$$353 \quad (4.7) \quad \|\nu\|_{\ell_2} := \left(\sum_{i=1}^{N-1} (z_{i+1} - z_i) F_\nu(z_i)^2 + F_\nu(z_N)^2 \right)^{1/2}$$

354 for all $\nu \in \mathcal{L}_z$. It is not difficult to see that $\|\cdot\|_{\ell_2}$ is a norm on \mathcal{L}_z and induces the
 355 metric ℓ_2 on \mathcal{P}_z . By taking the supremum over all state-action pairs, this property
 356 extends to $\bar{\ell}_2$.

357 Further we define the norm

$$358 \quad \|\nu\|_{\bar{\ell}_\infty} := \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \|\nu\|_{\ell_\infty} := \sup_{(x,a) \in \mathcal{X} \times \mathcal{A}} \max_{1 \leq i \leq N} |F_\nu(z_i)|$$

359 for all $\nu \in \mathcal{L}_z$. The inequalities

$$360 \quad (4.8) \quad \ell_2(\mu, \nu) = \|\mu - \nu\|_{\ell_2} \leq \sqrt{2V_{\max}} \|\mu - \nu\|_{\ell_\infty} \leq \sqrt{2V_{\max}}$$

361 hold for all μ and $\nu \in \mathcal{P}_z$. Lastly, since $\epsilon_k^{(x,a)}$ is the difference of a random probability
 362 measure and a probability measure in \mathcal{P}_z (see proof of Lemma 4.4), $F_{\epsilon_k^{(x,a)}}(z_N) = 0$
 363 holds, and thus $F_{E_k^{(x,a)}}(z_N) = 0$ also. The inequality

$$364 \quad (4.9) \quad \|E_k\|_{\bar{\ell}_2} \leq \sqrt{2V_{\max}} \|E_k\|_{\bar{\ell}_\infty}$$

365 follows from (4.7).

366 LEMMA 4.6. *For all $k \geq 1$, the inequalities*

$$367 \quad \|\eta_C - \eta_k\|_{\bar{\ell}_2} \leq \frac{\sqrt{\gamma\bar{\beta}}}{k} \sqrt{2V_{\max}} + \frac{1}{k} \sum_{j=1}^k \sqrt{\gamma}^{k-j} \|E_{j-1}\|_{\bar{\ell}_2}$$

368 *hold.*

369 *Proof.* Again, this is proved by induction. We use the fact that $\Pi_C \mathcal{T}^\pi$ is a $\sqrt{\gamma}$ -
 370 contraction in $\bar{\ell}_2$, substitute the equality from Lemma 4.5, and apply the norm in-
 371 equality (4.8).

372 For $k = 1$, the inequality holds as

$$\begin{aligned} 373 \quad \|\eta_C - \eta_1\|_{\bar{\ell}_2} &= \|\Pi_C \mathcal{T}^\pi \eta_C - \Pi_C \mathcal{T}^\pi \eta_0 + E_0\|_{\bar{\ell}_2} \\ 374 \quad &\leq \sqrt{\gamma} \|\eta_C - \eta_0\|_{\bar{\ell}_2} + \|E_0\|_{\bar{\ell}_2} \\ 375 \quad &\leq \sqrt{\gamma} \sqrt{2V_{\max}} + \|E_0\|_{\bar{\ell}_2} \\ 376 \quad &\leq \sqrt{\gamma\bar{\beta}} \sqrt{2V_{\max}} + \|E_0\|_{\bar{\ell}_2}. \end{aligned}$$

378 Assume that the equation holds for $k \geq 1$. It also holds for $k + 1$, since

$$\begin{aligned}
379 \quad & \|\eta_C - \eta_{k+1}\|_{\bar{\ell}_2} \\
380 \quad &= \left\| \Pi_C \mathcal{T}^\pi \eta_C - \frac{1}{k+1} (\Pi_C \mathcal{T}^\pi \eta_0 + k \Pi_C \mathcal{T}^\pi \eta_k - E_k) \right\|_{\bar{\ell}_2} \\
381 \quad &= \left\| \frac{1}{k+1} (\Pi_C \mathcal{T}^\pi \eta_C - \Pi_C \mathcal{T}^\pi \eta_0) + \frac{k}{k+1} (\Pi_C \mathcal{T}^\pi \eta_C - \Pi_C \mathcal{T}^\pi \eta_k) + \frac{1}{k+1} E_k \right\|_{\bar{\ell}_2} \\
382 \quad &\leq \frac{\sqrt{\gamma}}{k+1} \|\eta_C - \eta_0\|_{\bar{\ell}_2} + \frac{k\sqrt{\gamma}}{k+1} \|\eta_C - \eta_k\|_{\bar{\ell}_2} + \frac{1}{k+1} \|E_k\|_{\bar{\ell}_2} \\
383 \quad &\leq \frac{\sqrt{\gamma}}{k+1} \sqrt{2V_{\max}} + \frac{k\sqrt{\gamma}}{k+1} \left[\frac{\sqrt{\gamma}\bar{\beta}}{k} \sqrt{2V_{\max}} + \frac{1}{k} \sum_{j=1}^k \sqrt{\gamma}^{k-j} \|E_{j-1}\|_{\bar{\ell}_2} \right] + \frac{1}{k+1} \|E_k\|_{\bar{\ell}_2} \\
384 \quad &= \frac{\frac{\sqrt{\gamma}-\sqrt{\gamma^2}}{1-\sqrt{\gamma}}}{k+1} \sqrt{2V_{\max}} + \frac{\sqrt{\gamma^2}\bar{\beta}}{k+1} \sqrt{2V_{\max}} + \frac{1}{k+1} \sum_{j=1}^{k+1} \sqrt{\gamma}^{k+1-j} \|E_{j-1}\|_{\bar{\ell}_2} \\
385 \quad &= \frac{\sqrt{\gamma}\bar{\beta}}{k+1} \sqrt{2V_{\max}} + \frac{1}{k+1} \sum_{j=1}^{k+1} \sqrt{\gamma}^{k+1-j} \|E_{j-1}\|_{\bar{\ell}_2}, \\
386 \quad &
\end{aligned}$$

387 which concludes the proof. \square

388 **4.4. Step 4: Bounding the Error in Probability.** Applying the Hoeffding-
389 Azuma inequality is the crucial step in proving Theorem 3.1.

390 LEMMA 4.7 (Maximal Hoeffding-Azuma Inequality [5]). *Let $\mathcal{V} := \{V_1, \dots, V_T\}$ be
391 a martingale difference w.r.t. to the filtration \mathcal{F}_k ($\mathbb{E}[V_k | \mathcal{F}_{k-1}] = 0$) such that \mathcal{V} is
392 uniformly bounded by $L > 0$. Then for any $\epsilon > 0$, the inequality*

$$393 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} \left| \sum_{i=1}^k V_i \right| > \epsilon \right) \leq 2 \exp \left(\frac{-\epsilon^2}{2TL^2} \right)$$

394 holds.

395 LEMMA 4.8. *For all $\epsilon > 0$ and all time steps T , the inequality*

$$396 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} \|E_{k-1}\|_{\bar{\ell}_\infty} > \epsilon \right) \leq 2nN \exp \left(\frac{-\epsilon^2}{2T} \right)$$

397 holds.

398 *Proof.* Fix $(x, a) \in \mathcal{X} \times \mathcal{A}$ and define

$$400 \quad E_k^i := F_{E_k^{(x,a)}}(z_i) = \sum_{j=0}^k F_{\epsilon_j^{(x,a)}}(z_i).$$

401 By Lemma 4.4, $V_j = F_{\epsilon_j^{(x,a)}}(z_i)$, $j = 0, \dots, T$, is a martingale difference se-
402 quence w.r.t. \mathcal{F}_j and uniformly bounded by 1. Therefore, we can apply the maximal
403 Hoeffding-Azuma inequality, which takes the form

$$404 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} |E_{k-1}^i| > \epsilon \right) \leq 2 \exp \left(\frac{-\epsilon^2}{2T} \right).$$

405 By taking the union over all atoms, we have

$$\begin{aligned}
406 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} \left\| E_{k-1}^{(x,a)} \right\|_{\ell_\infty} > \epsilon \right) &= \mathbb{P} \left(\max_{1 \leq k \leq T} \max_{1 \leq i \leq N} |E_{k-1}^i| > \epsilon \right) \\
407 &= \mathbb{P} \left(\bigcup_{i=1}^N \left\{ \max_{1 \leq k \leq T} |E_{k-1}^i| > \epsilon \right\} \right) \\
408 &\leq 2N \exp \left(\frac{-\epsilon^2}{2T} \right). \\
409
\end{aligned}$$

410 Similarly, taking the union over all $(x, a) \in \mathcal{X} \times \mathcal{A}$, we find

$$\begin{aligned}
411 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} \|E_{k-1}\|_{\bar{\ell}_\infty} > \epsilon \right) &\leq 2nN \exp \left(\frac{-\epsilon^2}{2T} \right), \\
412
\end{aligned}$$

413 which concludes the proof. \square

414 **4.5. Step 5: Concluding the Proof of Theorem 3.1.**

415 *Proof of Theorem 3.1.* By Lemma 4.6 and inequality (4.9), we find

$$\begin{aligned}
416 \quad \|\eta_C - \eta_T\|_{\bar{\ell}_2} &\leq \frac{\sqrt{\gamma}\bar{\beta}}{T} \sqrt{2V_{\max}} + \frac{1}{T} \sum_{k=1}^T \sqrt{\gamma}^{T-k} \|E_{k-1}\|_{\bar{\ell}_2} \\
417 &\leq \frac{\sqrt{\gamma}\bar{\beta}}{T} \sqrt{2V_{\max}} + \frac{\bar{\beta}}{T} \sqrt{2V_{\max}} \max_{1 \leq k \leq T} \|E_{k-1}\|_{\bar{\ell}_\infty}. \\
418
\end{aligned}$$

419 By Lemma 4.8 the inequality

$$\begin{aligned}
420 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} \|E_{k-1}\|_{\bar{\ell}_\infty} > \epsilon \right) &\leq 2nN \exp \left(\frac{-\epsilon^2}{2T} \right) =: \delta
\end{aligned}$$

421 holds. Setting δ as above and solving for ϵ yields

$$\begin{aligned}
422 \quad \mathbb{P} \left(\max_{1 \leq k \leq T} \|E_{k-1}\|_{\bar{\ell}_\infty} \leq \sqrt{2T \log \frac{2nN}{\delta}} \right) &\geq 1 - \delta.
\end{aligned}$$

423 Therefore, with probability at least $1 - \delta$, we have

$$\begin{aligned}
424 \quad \bar{\ell}_2(\eta_C, \eta_T) = \|\eta_C - \eta_T\|_{\bar{\ell}_2} &\leq \sqrt{2V_{\max}}\bar{\beta} \left(\frac{\sqrt{\gamma}}{T} + \sqrt{\frac{2 \log \frac{2nN}{\delta}}{T}} \right),
\end{aligned}$$

425 which concludes the proof of the theorem. \square

426 *Proof of Corollary 3.2.* Define

$$\begin{aligned}
427 \quad T &:= \frac{C\bar{\beta}^2 V_{\max} \log \frac{2nN}{\delta}}{\epsilon^2}, \\
428 \quad t &:= \frac{\bar{\beta}^2 V_{\max} \log \frac{2nN}{\delta}}{\epsilon^2} \geq 1, \\
429
\end{aligned}$$

430 implying

$$\begin{aligned}
431 \quad \frac{1}{t} &\leq \frac{1}{\sqrt{t}}.
\end{aligned}$$

432 For $C = 2 + \sqrt{2} + 2\sqrt{1 + \sqrt{2}} \leq 6.53$, it follows that

$$433 \quad \bar{\ell}_2(\eta_C, \eta_T) \leq \epsilon\sqrt{2} \left(\frac{\sqrt{\gamma}}{C\sqrt{\log \frac{2nN}{\delta}}} + \sqrt{\frac{2}{C}} \right) \leq \epsilon\sqrt{2} \left(\frac{1}{C} + \sqrt{\frac{2}{C}} \right) \leq \epsilon. \quad \square$$

434
435 *Proof of Corollary 3.3.* After rearranging, we have

$$436 \quad \mathbb{P}(\bar{\ell}_2(\eta_C, \eta_T) > \epsilon) \leq 2nN \exp \left(\frac{\sqrt{\gamma}\epsilon}{\sqrt{2V_{\max}}\beta} - \frac{\gamma}{2T} - \frac{T\epsilon^2}{4V_{\max}\beta^2} \right).$$

437 As $\frac{\gamma}{2T} \geq 0$, we can omit this term. Since $\exp \left(-\frac{\epsilon^2}{4V_{\max}\beta^2} \right) < 1$, we find an inequality
438 of the form

$$439 \quad \mathbb{P}(\bar{\ell}_2(\eta_C, \eta_T) > \epsilon) \leq Cq^T, \quad C > 0, \quad 0 < q < 1.$$

440 Therefore $\sum_{T=0}^{\infty} \mathbb{P}(\bar{\ell}_2(\eta_C, \eta_T) > \epsilon) < \infty$, and by the Lemma of Borel-Cantelli we have
441 almost sure convergence. \square

442 **5. Policy Control.** Unfortunately, the analysis cannot be easily extended to
443 categorical distributions in the control case. There are several reasons.

444 First, the Bellman optimality operator \mathcal{T} is not a contraction in $\bar{\ell}_2$. Bellemare et
445 al. [1] provided a counter example for the Wasserstein distance that also works for
446 the Cramér distance. Therefore Lemma 4.6 does not hold if \mathcal{T}^π is replaced by \mathcal{T} .

447 Nevertheless, we consider the update rule for the control case,

$$448 \quad (5.1) \quad \eta_{k+1} := \eta_k + \alpha_k(\Pi_C \mathcal{T}_k^{\pi_{k-1}} \eta_{k-1} - \eta_k) + (1 - \alpha_k)(\Pi_C \mathcal{T}_k^{\pi_k} \eta_k - \Pi_C \mathcal{T}_k^{\pi_{k-1}} \eta_{k-1}).$$

449 Here π_k denotes the greedy policy with respect to the expected values of η_k and it
450 holds that $\mathcal{T}_k^{\pi_k} \eta_k = \mathcal{T}_k \eta_k$.

451 It can be shown that such update rules produce the same expected values as their
452 value-based algorithmic counterpart [7]. Therefore, we can be sure that $Q_k(x, a) :=$
453 $\mathbb{E}_{Z \sim \eta_k^{(x,a)}} [Z]$ converges to the unique optimal value function Q^* , because Q_k satisfies
454 equation (2.8).

455 If we assume a unique optimal policy π^* , then Q_k comes close enough to Q^* such
456 that $\pi_k = \pi^*$ after some time and it remains to evaluate the return distributions of
457 π^* for which convergence holds. This is the reasoning Rowland et al. [8] used to prove
458 their control theorem.

459 This approach does not work in the present control case, as the update rule (5.1)
460 does not necessarily yield probability measures anymore, which can be seen as follows
461 by revisiting the proof of Lemma 4.3 and calculating

$$462 \quad \mathcal{D}_{k+1}[\eta_{k+1}, \eta_k]^{(x,a)}$$

$$463 \quad = (k+1)\Pi_C \mathcal{T}_{k+1}^{\pi_{k+1}} \eta_{k+1}^{(x,a)} - k\Pi_C \mathcal{T}_{k+1}^{\pi_k} \eta_k^{(x,a)}$$

$$464 \quad = (k+1)\Pi_C \mathcal{T}_{k+1}^{\pi_{k+1}} \left(\frac{k}{k+1} \eta_k + \frac{1}{k+1} \mathcal{D}_k[\eta_k, \eta_{k-1}] \right)^{(x,a)} - k\Pi_C \mathcal{T}_{k+1}^{\pi_k} \eta_k^{(x,a)}$$

$$465 \quad = k\Pi_C \mathcal{T}_{k+1}^{\pi_{k+1}} \eta_k^{(x,a)} + \Pi_C \mathcal{T}_{k+1}^{\pi_{k+1}} \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)} - k\Pi_C \mathcal{T}_{k+1}^{\pi_k} \eta_k^{(x,a)}$$

467 for the control case. But if $\pi_{k+1} \neq \pi_k$, this is not equal to $\Pi_C \mathcal{T}_{k+1}^{\pi_{k+1}} \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)}$
468 in general and hence, $\frac{k}{k+1} \eta_k^{(x,a)} + \frac{1}{k+1} \mathcal{D}_k[\eta_k, \eta_{k-1}]^{(x,a)}$ is not necessarily an element
469 of $\mathcal{P}(\mathcal{P}_z)$, meaning that we now obtain signed measures in the general case.

470 In order to fix this problem, one could alter the update rule to become

471 (5.2) $\eta_{k+1} := \eta_k + \alpha_k(\Pi_C \mathcal{T}_k^{\pi_k} \eta_{k-1} - \eta_k) + (1 - \alpha_k)(\Pi_C \mathcal{T}_k^{\pi_k} \eta_k - \Pi_C \mathcal{T}_k^{\pi_k} \eta_{k-1}).$

472 With this changed definition, Lemma 4.3 holds again, but we run into different prob-
 473 lems. The first problem is that Lemma 4.5 does not hold any more, as we now have
 474

475
$$\eta_k = \frac{1}{k}(\Pi_C \mathcal{T}^{\pi_0} \eta_0 + (k-1)\Pi_C \mathcal{T}^{\pi_{k-1}} \eta_{k-1} - E_{k-1})$$

 476
$$+ \frac{1}{k} \sum_{j=0}^{k-1} (j-1)(\Pi_C \mathcal{T}^{\pi_{j-1}} \eta_{j-1} - \Pi_C \mathcal{T}^{\pi_j} \eta_{j-1}).$$

 477

478 But if the policies do not change anymore after time step T , the summands are zero for
 479 $k > T$ and the second term becomes small as k tends to infinity. Thus, $\eta_k \approx \Pi_C \mathcal{T} \eta_{k-1}$
 480 holds again, which indicates this problem to be minor. However, it leads to the second,
 481 more serious, problem, namely showing that the expected values of $\eta_k^{(x,a)}$ obtained by
 482 (5.2) still converge to the optimal value function Q^* .

483 Nevertheless, this adjusted update rule shows good experimental results, as dis-
 484 cussed in the next section.

485 **6. Experimental Results.**

486 **6.1. Combination-Lock.** Consider the combination-lock environment [4]. We
 487 have a set of 500 states x_i , which are arranged in a chain. In each state, we can choose
 488 between two actions **left** or **right**, see Figure 1. Choosing **right** takes the agent to
 489 state x_{i+1} , but yields a reward of -0.01 . Taking **left** takes the agent to a previous
 490 state with probability $p(x_k|x_i, \text{left}) \propto \frac{1}{i-k}$ and yields reward 0. Transitioning to the
 491 goal state x_{500} gives the reward $+15$.

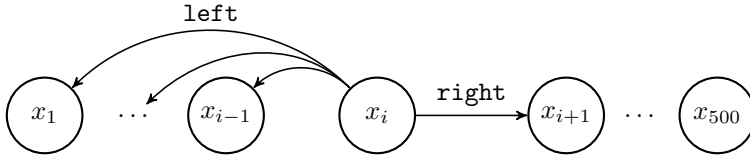


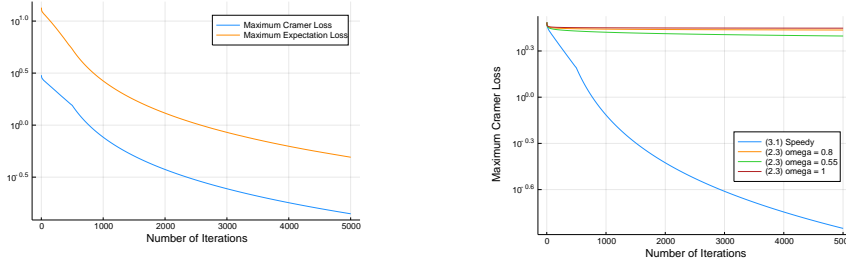
Fig. 1: Combination-lock environment

492 The action **right** brings us closer to the goal state, but yields a negative reward,
 493 whereas the action **left** has no immediate negative reward, but moves us further
 494 from x_{500} . The rewards are set up such that choosing **right** in all states is the unique
 495 optimal policy. This makes an interesting control problem, because the long chain has
 496 to be essentially solved right to left. It is also a good benchmark for policy evaluation,
 497 because the trajectories are long and when choosing **left** there are a lot of possible
 498 successor states.

499 In the experiment, $\gamma = 0.999$ and 51 equally spaced atoms or grid points between
 500 -10 and 15 were chosen. The SCPE algorithm was run 10 times for 5000 iterations
 501 with random initial distributions (51 random numbers were drawn independently from
 502 the uniform distribution $[0, 1]$ for all (x, a) and then divided by their sum to form
 503 probabilities). For comparison, the TD update rule (2.3) with polynomial learning
 504 rates $\omega \in \{0.55, 0.8, 1\}$ was tested. The limiting return distribution η_C was estimated
 505 by performing SCPE for 50 000 iterations, denoted by $\hat{\eta}_C$.

506 In Figure 2a, the maximum Cramér distance $\bar{\ell}_2(\eta_k, \hat{\eta}_C)$ to the estimated limiting
 507 return distribution function and the maximum absolute distance of the corresponding
 508 expected returns, averaged over the 10 runs, are shown. This confirms that indeed
 509 about the same sample complexity holds in both cases.

510 In Figure 2b the clear performance benefit of the speedy update rule (3.1) over
 511 the TD one (2.3) is visible. This plot resembles the results of Ghavamzadeh et al. [4].

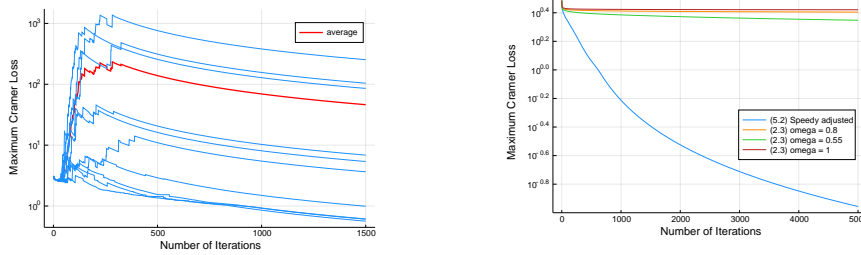


(a) Convergence in Cramér distance vs. convergence in expectation

(b) SCPE and polynomial learning rates

Fig. 2: Policy evaluation in the combination-lock environment.

512 The same experiment was performed in the control case. The instability problem
 513 of using the unadjusted update rule (5.1) is illustrated in Figure 2a. Here, measures
 514 with negative probabilities were indeed produced. Using the adjusted update rule (5.2)
 515 yields almost exactly the same performance improvements as in the policy evaluation
 516 case, see Figures 2b and 3b.



(a) Instability of unchanged algorithm

(b) Improved convergence for adjusted algorithm

Fig. 3: Q-learning in the combination-lock environment.

517 **6.2. Gridworld.** In order to put the adjusted update rule (5.2) to the test, we
 518 investigated the convergence of the expected values to the optimal value function Q^*
 519 in an environment with multiple optimal policies. We consider an $n \times n$ gridworld,
 520 where the agent can move up, down, left, and right. If the agent moves to the cell with
 521 coordinates $(x, y) \in \{1, \dots, n\}^2$, it receives reward $\pm(|(n-x+1)-y|+1)$ with equal
 522 probability. Only at the goal cell (n, n) , the positive reward n is always obtained.
 523 Figure 4 shows an overview of this environment.

524 The difficulty for the agent is to recognise that wandering around in the environ-
 525 nment gives an expected return of zero and that the optimal strategy is to reach the
 526 goal cell as quickly as possible. This can be done along multiple paths in the grid and
 527 lucky immediate rewards causes the agent to often change direction.

$\pm n$	$\pm(n-1)$...	± 3	± 2	± 1
$\pm(n-1)$	$\pm(n-2)$...	± 2	± 1	± 2
\vdots	\vdots	\ddots	\vdots	\vdots	\vdots
± 2	± 1	...	$\pm(n-3)$	$\pm(n-2)$	$\pm(n-1)$
± 1	± 2	...	$\pm(n-2)$	$\pm(n-1)$	n

Fig. 4: Gridworld with rewards given at each cell.

528 For this environment we used the same experiment setup as in Section 6.1 with
 529 $n = 25$ and $\gamma = 0.9$. While it was possible to compute the expected values simply
 530 from the categorical distributions for update rules (2.3) and (5.2), this was not the
 531 case for rule (5.1) in the gridworld environment. The unadjusted update rule (5.1)
 532 led to such instabilities that the signed probabilities under- and overflowed the 64 bit
 533 double value range. For this reason, we directly used the update rule for Q values
 534 (2.8) in this case, where the initial Q -tables were uniformly sampled from $[-n, n]$.

535 In Figure 5, the maximum absolute difference to the theoretical optimal value
 536 function Q^* is shown for both environments. For the combination-lock environment,
 537 the expected values were obtained from the distributional updates of the control
 538 experiment of Section 2; for the gridworld, the expected values were obtained as
 539 described in the last paragraph. While the adjusted update rule is slightly better
 540 than the unadjusted one in the combination-lock environment, it lags behind in the
 541 gridworld example.

542 Both the adjusted (5.2) and the unadjusted (5.1) speedy Q-learning update rules,
 543 are convincingly faster than the standard Q-learning update rule in both environ-
 544 nments. Further, the greedy policies changed up to around time step 600 in the
 545 combination-lock environment, whereas in the gridworld they change in over 85%
 546 of the steps up to the last iteration. This suggests that the adjusted update rule is
 547 still robust under frequent policy changes, but they may slow down convergence.

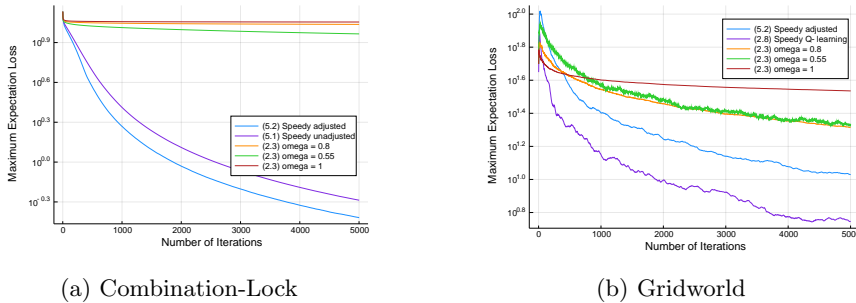


Fig. 5: Comparison of the adjusted (5.2) versus the unadjusted (5.1) update rule by the maximum absolute difference to the theoretical optimal value function.

548 **7. Conclusions.** In this paper, speedy Q-learning was extended from the value-
 549 based case to categorical distributions. For evaluating a fixed policy, PAC bounds
 550 in terms of the Cramér distance were established. This led to the observation that
 551 even though the computational and space complexity scale linearly in the number of
 552 atoms N , the sample complexity scales only logarithmically in N . Thus, switching
 553 from standard RL to distributional RL or increasing the accuracy of the distribution
 554 approximation yields only a small penalty in terms of transition samples needed. An
 555 application in two simple environments confirmed the theoretical results empirically.

556 The reasons as to why the finite-time analysis could not be easily extended to
 557 the case of policy control were stated. Experiments showed that a slight modification
 558 to the update rule results in the same performance improvements as in the policy
 559 evaluation case. An in-depth analysis of this adjusted updated rule remains for future
 560 work.

561

REFERENCES

- 562 [1] M. G. BELLEMARE, W. DABNEY, AND R. MUNOS, *A distributional perspective on reinforcement*
 563 *learning*, in Proceedings of the 34th International Conference on Machine Learning
 564 – Volume 70, ICML 2017, JMLR.org, 2017, p. 449–458.
- 565 [2] R. BELLMAN, *Dynamic Programming*, Princeton University Press, 1957.
- 566 [3] E. EVEN-DAR AND Y. MANSOUR, *Learning rates for Q-learning*, J. Mach. Learn. Res., 5 (2004),
 567 p. 1–25.
- 568 [4] M. GHAVAMZADEH, H. J. KAPPEN, M. G. AZAR, AND R. MUNOS, *Reinforcement learning with*
 569 *a near optimal rate of convergence*, technical report, INRIA, Oct. 2011. ID 00636615v2.
- 570 [5] M. GHAVAMZADEH, H. J. KAPPEN, M. G. AZAR, AND R. MUNOS, *Speedy Q-learning*, in Advances
 571 in Neural Information Processing Systems 24, J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett,
 572 F. Pereira, and K. Q. Weinberger, eds., Curran Associates, Inc., 2011, pp. 2411–2419.
- 573 [6] O. KALLENBERG, *Random Measures, Theory and Applications*, Springer, 2017.
- 574 [7] C. LYLE, P. S. CASTRO, AND M. G. BELLEMARE, *A comparative analysis of expected and*
 575 *distributional reinforcement learning*, CoRR, abs/1901.11084 (2019).
- 576 [8] M. ROWLAND, M. BELLEMARE, W. DABNEY, R. MUNOS, AND Y. W. TEH, *An analysis of*
 577 *categorical distributional reinforcement learning*, in Proceedings of the Twenty-First Inter-
 578 national Conference on Artificial Intelligence and Statistics, A. Storkey and F. Perez-Cruz,
 579 eds., vol. 84 of Proceedings of Machine Learning Research, Playa Blanca, Lanzarote, Ca-
 580 nary Islands, 9–11 Apr 2018, PMLR, pp. 29–37.
- 581 [9] R. SUTTON, *Learning to predict by the method of temporal differences*, Machine Learning, 3
 582 (1988), pp. 9–44.
- 583 [10] R. S. SUTTON AND A. G. BARTO, *Reinforcement Learning: An Introduction*, The MIT Press,
 584 second ed., 2018.
- 585 [11] C. J. C. H. WATKINS AND P. DAYAN, *Q-learning*, Machine Learning, (1992), pp. 279–292.