

Effective Exploration for Deep Reinforcement Learning via Bootstrapped Q-Ensembles under Tsallis Entropy Regularization

Gang Chen^{*1}, Yiming Peng^{†1}, and Mengjie Zhang^{‡1}

¹School of Engineering and Computer Science,
Victoria University of Wellington, New Zealand

September 6, 2019

Abstract

Recently *deep reinforcement learning* (DRL) has achieved outstanding success on solving many difficult and large-scale RL problems. However the high sample cost required for effective learning often makes DRL unaffordable in resource-limited applications. With the aim of improving sample efficiency and learning performance, we will develop a new DRL algorithm in this paper that seamlessly integrates *entropy-induced* and *bootstrap-induced* techniques for efficient and deep exploration of the learning environment. Specifically, a general form of *Tsallis entropy regularizer* will be utilized to drive entropy-induced exploration based on efficient approximation of optimal action-selection policies. Different from many existing works that rely on action dithering strategies for exploration, our algorithm is efficient in exploring actions with clear exploration value. Meanwhile, by employing an ensemble of Q-networks under varied Tsallis entropy regularization, the diversity of the ensemble can be further enhanced to enable effective bootstrap-induced exploration. Experiments on Atari game playing tasks clearly demonstrate that our new algorithm can achieve more efficient and effective exploration for DRL, in comparison to recently proposed exploration methods including Bootstrapped Deep Q-Network and UCB Q-Ensemble.

Keywords— Reinforcement Learning, Deep Learning, Q-Ensemble, Tsallis Entropy

1 Introduction

In recent years, *deep reinforcement learning* (DRL) has been extensively and successfully utilized by computer systems to autonomously learn to solve many challenging problems such as

^{*}aaron.chen@ecs.vuw.ac.nz

[†]yiming.peng@ecs.vuw.ac.nz

[‡]mengjie.zhang@ecs.vuw.ac.nz

robotics control [LHP⁺15, SLA⁺15], video game playing [MKS⁺15, HGS16], and road traffic management [LLW16]. However, in order to achieve its learning goals, an RL agent must often use a huge amount of sampled data to train its *deep neural networks* (DNNs). Since data sampling is realized through direct trial-and-error interactions with the learning environment, the high *sample cost* usually makes DRL unaffordable in resource-limited applications [WMG⁺17, WBH⁺16, CPZ18a].

In order to improve sample efficiency, an RL agent must carefully manage its exploration of the learning environment. Osband *et al.* recently proposed the idea of “*deep exploration*” to emphasize on the requirement for the agent to learn effectively within a reasonable time frame by considering not only the immediate benefits of taking any action but also the long-term impact of the action on future learning, thereby properly synthesizing efficient exploration with effective generalization [OR16, ORW14, ORWR17]. Guided by this requirement, Bootstrapped Deep Q-Network (Bootstrapped DQN) has been proposed lately to drive deep and efficient exploration [OBPR16].

Bootstrapped DQN was inspired by the posterior sampling method for RL with near-optimal regret bounds [ORR13]. However, instead of sampling and solving numerous *Markov Decision Processes* (MDPs), Bootstrapped DQN approximates a posterior model over optimal Q-functions (also known as the state-action value functions) at much affordable computation cost. This is shown to easily outperform *action dithering strategies* for exploration such as ϵ -greedy or softmax action sampling techniques [Kak03, Str07]. For this purpose, an ensemble of randomly initialized Q-networks (or Q-functions) will be maintained consistently during RL. Empirical results showed that effective deep exploration can be achieved in practice by randomly choosing one of the Q-networks to guide multi-step interactions with the learning environment.

Besides Bootstrapped DQN, the UCB Q-Ensemble method proposed in [CSAS17] also relies on learning concurrently an ensemble of Q-networks. However it adopts an approximated *upper-confidence bound* over Q-values produced by these Q-networks to steer exploration. Although highly competitive performance has been witnessed on Atari game playing tasks, theoretical studies suggest that precise calculation of such confidence bounds can be computationally intractable [RR13, RR14].

Similar to Bootstrapped DQN and UCB Q-Ensemble, we employ an ensemble of Q-networks to achieve deep exploration. However, without relying on actions with either the highest Q-values or upper-confidence bounds for exploration, we generalize action selection by studying policies under entropy regularization. This generalization enables us to develop a new form of optimal stochastic policies, thereby relieving the dependency on randomly initialized Q-networks as the main source of randomness for deep exploration [CSAS17, OBPR16].

In the literature, Shannon entropy is frequently utilized to regularize action selection, giving rise to optimal policies that exhibit *softmax action-selection behaviors* [OMKM17, NNXS17, HTAL17, SAC17]. While softmax distributions naturally bring stochasticity to deep exploration, they are prone to assigning non-negligible probability mass to actions with negligible exploration value [LCO17, NCG18].

Tsallis entropy is an important extension of Shannon entropy [PP93, Tsa94]. A special case of *Tsallis entropy* has been studied in [LCO17] to tackle sparse MDP problems. When applied to DRL, Tsallis entropy allows an RL agent to concentrate on exploring actions that deserve further exploration. Due to this reason, *Tsallis entropy regularization is deemed a key mechanism for efficient deep exploration* in this paper. Moreover, without being restricted to any specific setting of Tsallis entropy as in [LCO17, NCG18], a general form of Tsallis entropy will be studied the first time in literature to guide DRL.

Based on computationally efficient approximation of optimal policies under general Tsallis entropy regularization, a new deep exploration algorithm involving an ensemble of deep Q-networks will be further developed in this paper. The newly proposed algorithm will be called

the *Bootstrapped Q-Ensemble under Tsallis Entropy Regularization* (BQETR) algorithm. Each Q-network in BQETR adopts a different setting of the Tsallis entropy regularizer in order to achieve high ensemble diversity. Meanwhile, the *regularization coefficient* is kept the same for all Q-networks and will be gradually reduced to 0 as an RL agent gains increasingly more experience from its learning environment. In this way we expect to seamlessly integrate *entropy-induced exploration* with *bootstrap-induced exploration* for effective RL.

Empirical studies have been performed on benchmark Atari game playing tasks. Our experiment results clearly show that BQETR has achieved significantly better sample efficiency and performance than Bootstrapped DQN and UCB Q-Ensemble. We therefore believe that BQETR is an effective and efficient method for deep exploration and RL.

2 Related Works

Huge efforts have been devoted to developing efficient and effective exploration strategies for RL. Particularly, provably efficient exploration techniques have been studied based on the idea of Bayes optimal policies [GMP⁺15] and clearly revealed the importance of multi-step exploration [KS02]. Further studies along this line also demonstrated the inefficiency of ϵ -greedy and softmax exploration techniques on large RL problems [BT02, SLW⁺06, AO07, DB15]. In view of the fact that many existing DRL algorithms rely on such simple methods for exploration [MKS⁺15, HGS16], developing new exploration methods for effective DRL has great value both in theory and in practice.

Among all the exploration techniques proposed so far, a notable series of research works clearly highlighted the advantages of exploration through randomized value functions [ORR13, ORW14, OR16, ORWR17]. Specifically, the randomized least-squares value iteration (RLSVI) algorithms proposed in [ORWR17] extended traditional least-squares value iteration methods through randomly sampling *statistically plausible value functions*. However, efficient sampling often requires value functions to be linear with respect to their parameters and may not be suitable for DRL [ORW14, ORWR17]. To cope with this issue, Bootstrapped DQN has been developed recently in [OBPR16] to approximately sample value functions modeled as DNNs.

Besides bootstrapping, to achieve entropy-induced exploration, the training of action-selection policies is often reshaped by a Shannon entropy regularizer, resulting in various soft-Q style algorithms for DRL [NNXS17, OMKM17, HTAL17]. For example, Nachum *et al.* developed an off-policy algorithm based on a multi-step consistency equation for entropy-regularized RL [NNXS17]. Haarnoja *et al.* conducted research on soft-Q learning in high-dimensional action spaces [HTAL17]. In [OMKM17, SAC17], policy gradient training is shown as equivalent to soft-Q learning. This insightful understanding enables an RL agent to combine policy gradient with Q-learning for effective sample reuse [OMKM17]. Different from these research works, efficient entropy-induced exploration is realized in this paper through Tsallis entropy regularization [LCO17, NCG18].

This paper is similar to [LCO17, NCG18] since they all leverage on Tsallis entropy to regularize policy optimization. However, different from these research works that studied only a specific setting of Tsallis entropy, we will consider general forms of Tsallis entropy so as to maintain strong diversity in a Q-ensemble. It also allows us to treat Shannon entropy regularization as a special case of our research. Moreover, the integrated use of Tsallis entropy and bootstrapping mechanism for deep exploration further separates this paper apart from most of the previous works.

3 Entropy Regularized Q-Learning

In this section, we will introduce the RL problem first, followed by a quick review of Q-learning and policy gradient learning techniques. Afterwards, a new deep Q-learning algorithm under Tsallis entropy regularization will be developed. The *Bellman residue* of the newly proposed algorithm will also be analyzed under an extreme circumstance.

3.1 The Reinforcement Learning Problem

This paper studies general RL problems that can be described by an MDP with an arbitrary set of states $s \in \mathbb{S}$ and a finite set of actions $a \in \mathbb{A}$ [OMKM17]. Such problems appear frequently in literature including robotics control and video game playing [MKS⁺15, HGS16, CPZ18b]. At each time step t , an RL agent observes its environment and determines its current state s_t . It subsequently selects and performs an action a_t , driving the environment to move to its next state s_{t+1} with a probability $P(s_t, a_t, s_{t+1})$ which is unknown to the agent. Meanwhile, a scalar reward $r(s_t, a_t)$ is provided as the immediate feedback of performing action a_t . Starting from any initial state s_0 , the agent is required to perform a long (sometimes infinite) sequence of actions in order to obtain the *discounted total return* defined below

$$J(\pi) = \mathbb{E}_{s_0, \pi} \left(\sum_{t=0}^{\infty} \gamma^t r_t \right) \quad (1)$$

where the expectation in (1) is conditional on initial state s_0 and π . Here π refers to a *stochastic action-selection policy* that the RL agent follows to determine the probability $\pi(s_t, a)$ of performing any action a in any state s_t . Obviously $\sum_{a \in \mathbb{A}} \pi(s_t, a) = 1$ is an important constraint for π to be well-defined. Moreover, γ takes its value in $[0, 1)$ and serves as a discount factor for the RHS of (1) to be meaningful. With an MDP described as above, the ultimate goal of RL is hence to identify the *optimal policy* π^* that maximizes $J(\pi^*)$, i.e.

$$\pi^* = \arg \max_{\pi} J(\pi) \quad (2)$$

In an effort to learn π^* , an RL agent may choose to first learn the Q-function with respect to some non-optimal policy π , as defined below

$$Q^\pi(s_t, a_t) = \mathbb{E}_{\pi, s_t, a_t} \left(\sum_{k=0}^{\infty} \gamma^k r_{k+t} \right) \quad (3)$$

Given Q^π in (3), the value of state s_t under policy π is determined further as

$$V^\pi(s_t) = \mathbb{E}_{s_t, \pi} \left(\sum_{k=0}^{\infty} \gamma^k r_{k+t} \right) = \sum_{a \in \mathbb{A}} \pi(s_t, a) Q^\pi(s_t, a) \quad (4)$$

To ease discussion, we denote the value of any state s under the optimal policy π^* as $V^*(s)$. Accordingly $Q^*(s, a)$ represents the maximum Q-value achievable as a result of performing action a in state s .

3.2 Q-Learning and Policy Gradient

In value function based methods for DRL, the Q-function (or V-function) is represented as a DNN with numerous parameters. Through updating these parameters, we can bring the Q-value

outputs from such deep Q-networks as close to the fixed point of the *Bellman equation* as possible. Two versions of the Bellman equation are typically studied in the literature. Each version is associated with a different *Bellman operator*, i.e. \mathcal{T}^π and \mathcal{T}^* , as defined below.

$$\mathcal{T}^\pi Q^\pi(s, a) = r(s, a) + \gamma \sum_{s'} P(s, a, s') \sum_{b \in \mathbb{A}} \pi(s', b) Q^\pi(s', b) \quad (5)$$

$$\mathcal{T}^* Q^*(s, a) = r(s, a) + \gamma \sum_{s'} P(s, a, s') \max_{b \in \mathbb{A}} Q^*(s', b) \quad (6)$$

Clearly \mathcal{T}^π in (5) applies to Q^π which is the fixed point of the Bellman equation $\mathcal{T}^\pi Q(s, a) = Q(s, a)$. Likewise \mathcal{T}^* in (6) applies to Q^* which is the fixed point of the Bellman equation $\mathcal{T}^* Q(s, a) = Q(s, a)$. It is well-known that both \mathcal{T}^π and \mathcal{T}^* are γ -contraction mappings in the sup-norm and are suitable to drive value function learning [Ber95]. Specifically, in DQN [MKS⁺15], an approximation of \mathcal{T}^* based on a batch of previously sampled state transition data $\mathcal{B} = \{(s, a, s', r(s, a))\}$ is utilized to update the Q-network parameters in the direction of minimizing the loss function below

$$L^*(\theta) = \sum_{(s, a, s', r) \in \mathcal{B}} \left(Q_\theta(s, a) - \tilde{\mathcal{T}}^*(s, a, s', r) \right)^2 \quad (7)$$

To avoid overestimation in the original design of DQN, Double DQN is typically used in practice by employing a separate target Q-network to calculate $\tilde{\mathcal{T}}^*$ in (7) [HGS16]. The equation below shows more details.

$$\tilde{\mathcal{T}}^*(s, a, s', r) = r(s, a) + \gamma Q_{\theta^-} \left(s', \arg \max_{b \in \mathbb{A}} Q_\theta(s', b) \right) \quad (8)$$

where in (7) and (8), Q_θ refers to the Q-network parameterized by θ and Q_{θ^-} stands for the target Q-network parameterized by θ^- . Clearly by minimizing L^* in (7), Q_θ has the aim to approximate Q^* . Similar loss function has also been defined for Q_θ to precisely estimate Q^π with respect to arbitrary policy π .

Under the general actor-critic framework for RL, assume that a stochastic policy π is implemented as a DNN parameterized by ω . According to the policy gradient theorem [SMSM00], ω should be updated in the direction of

$$\frac{\partial J(\pi)}{\partial \omega} = \mathbb{E}_{(s, a) \sim \pi} \left(Q^\pi(s, a) \frac{\partial \log \pi(s, a)}{\partial \omega} \right) \quad (9)$$

Here the expectation is taken over all possible state-action pairs (s, a) with probability $d^\pi(s)\pi(s, a)$ and $d^\pi(s)$ gives the *discounted distribution of states* defined in [SMSM00].

3.3 Q-Learning under Tsallis Entropy Regularization

While training policy networks based on (9), in order to prevent a policy from converging too fast and therefore leaving no opportunity for future exploration, it is a common practice to introduce an extra *entropy regularizer*. As a consequence, the policy network parameters ω can be updated according to

$$\Delta \omega \propto \mathbb{E}_{(s, a) \sim \pi} \left(Q_\theta^\pi(s, a) \frac{\partial \log \pi(s, a)}{\partial \omega} + \alpha \frac{\partial H^\pi(s)}{\partial \omega} \right) \quad (10)$$

Algorithm 1 An Algorithm for Deep Q-Learning under Tsallis Entropy Regularization

- 1: **Input:** a Q-network, q value for the Tsallis entropy regularizer, α_0 for the initial regularization coefficient, and a replay buffer \mathcal{B} that stores past state-transition samples for training
 - 2: **for** each problem episode **do**:
 - 3: Obtain initial state s_0 from environment
 - 4: **for** $t = 1, \dots$ until end of episode **do**:
 - 5: Sample action a_t according to (27)
 - 6: Perform a_t
 - 7: Add (s_t, a_t, s_{t+1}, r_t) to \mathcal{B}
 - 8: **if** learning interval is reached **do**:
 - 9: Sample mini-batch from \mathcal{B}
 - 10: Update Q-network to minimize $L^{\pi_\alpha^*}(\theta)$ in the mini-batch
 - 11: Reduce α linearly by Δ_α until 0
-

where H^π denotes the entropy of policy π and $\alpha > 0$ is the *entropy regularization coefficient*. Previously, Shannon entropy as defined below is frequently utilized for regularized policy training.

$$H_I^\pi(s) = - \sum_{a \in \mathbb{A}} \pi(s, a) \log \pi(s, a) \quad (11)$$

Subject to entropy regularization in (11) with coefficient α in (10), it can be shown that the optimal policy π_α^* and the corresponding optimal Q-function $Q^{\pi_\alpha^*}$ obey the following equation [SAC17],

$$\pi_\alpha^*(s, a) = \frac{\exp(Q^{\pi_\alpha^*}(s, a)/\alpha)}{\sum_{b \in \mathbb{A}} \exp(Q^{\pi_\alpha^*}(s, b)/\alpha)} \quad (12)$$

In this paper, instead of using the softmax distribution in (12) to guide entropy-induced exploration which can exhibit poor efficiency in practice [ORWR17], we decide to study Tsallis entropy based regularizer as defined below.

$$H_q^\pi(s) = \frac{1}{q-1} \left(1 - \sum_{a \in \mathbb{A}} \pi(s, a)^q \right) \quad (13)$$

It can be shown that $\lim_{q \rightarrow 1} H_q^\pi(s) = H_I^\pi(s)$, for any $s \in \mathbb{S}$. We can hence consider soft-Q learning as demonstrated by (12) as a special case of our new Q-learning method. In [LCO17, NCG18], a specific setting of (13) with $q = 2$ has been studied to derive a fixed-form representation of the optimal policy. In this paper, on the other hand, we are interested in the general form of Tsallis entropy in (13) with $q > 1$.

Analogous to the analysis presented in [OMKM17], we can represent the RHS of (10) as $f(\omega)$. Meanwhile let $g^\pi(s) = \sum_{a \in \mathbb{A}} \pi(s, a)$. Clearly when ω in (10) reaches its fixed point (or optima), no further updating of ω in the direction of $f(\omega)$ is possible without violating the constraint that $g^\pi(s) = 1$ for any $s \in \mathbb{S}$. This means that, with the optimal policy parameters ω^* , $f(\omega^*)$ belongs to the span of the vectors $\{\frac{\partial g^\pi(s)}{\partial \omega}\}$, i.e.

$$f(\omega^*) = \sum_{s \in \mathbb{S}} \lambda(s) \frac{\partial g^\pi(s)}{\partial \omega} \Big|_{\omega=\omega^*} \quad (14)$$

where for every state s , the Lagrange multiplier $\lambda(s)$ in (14) ensures that $g^\pi(s) = 1$. Meanwhile we can determine $\frac{\partial H_q^\pi(s)}{\partial \omega}$ as

$$\begin{aligned} \frac{\partial H_q^\pi(s)}{\partial \omega} &= \frac{q}{q-1} \sum_{a \in \mathbb{A}} \pi(s, a)^{q-1} \frac{\partial \pi(s, a)}{\partial \omega} \\ &= \frac{q}{q-1} \sum_{a \in \mathbb{A}} \pi(s, a)^q \frac{\partial \log \pi(s, a)}{\partial \omega} \end{aligned} \quad (15)$$

By substituting (15) into (10) and also taking into account (14), the optimal condition for policy parameters ω^* becomes

$$\begin{aligned} &\mathbb{E}_{(s, a) \sim \pi} \left(\left(Q^\pi(s, a) - \frac{\alpha q}{q-1} \pi(s, a)^{q-1} - c(s) \right) \frac{\partial \log \pi(s, a)}{\partial \omega} \right) \\ &= 0 \end{aligned} \quad (16)$$

where $c(s)$ stands for $\lambda(s)$ in (14) adjusted according to the discounted distribution of state s . To solve the equation in (16), similar to [OMKM17], it is eligible to consider each state s separately. Particularly, in any state s and $\forall a \in \mathbb{A}$, we have

$$Q^\pi(s, a) - \frac{\alpha q}{q-1} \pi(s, a)^{q-1} - c(s) = 0 \text{ or } \frac{\partial \pi(s, a)}{\partial \omega} = 0 \quad (17)$$

It is straightforward to verify the solution below of (17),

$$\pi_\alpha^*(s, a) = \sqrt[q-1]{\max \left(\left(\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - \frac{c(s)}{\alpha} \right), 0 \right) \frac{q-1}{q}} \quad (18)$$

with π_α^* representing the optimal policy of the entropy-regularized policy gradient learning problem described in (10). $Q^{\pi_\alpha^*}$ stands for the respective Q-function for policy π_α^* . Meanwhile, $c(s)$ in (18) ensures that the condition $g^\pi(s) = 1$ holds consistently. Notice that for certain action a , it is possible for $Q^{\pi_\alpha^*}(s, a) - c(s) < 0$. For such an action, the validity of (18) is ensured by letting $\pi^*(s, a) = 0$, therefore $\frac{\partial \pi^*(s, a)}{\partial \omega} = 0$. In other words, only a portion of actions in \mathbb{A} may be explored in any state s , thereby encouraging efficient exploration. Particularly, when $q = 2$, it can be shown that [LCO17]

$$c(s) = \alpha \frac{\sum_{a \in S(s)} \frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - q}{\|S(s)\|} \quad (19)$$

with $S(s)$ representing the set of actions with non-zero chance of exploration in state s , as determined below.

$$S(s) = \left\{ a_i \mid q + i \frac{Q^{\pi_\alpha^*}(s, a_i)}{\alpha} > \sum_{j=1}^i \frac{Q^{\pi_\alpha^*}(s, a_j)}{\alpha} \right\} \quad (20)$$

where a_i denotes the action with the i -th highest Q-value in state s . If $q \neq 2$, closed-form representation of $c(s)$ and $S(s)$ may not exist. Therefore, in order to estimate π_α^* , we have developed two efficient approximation techniques in our Q-Learning algorithm. Specifically,

because whenever $\pi_\alpha^*(s, a) > 0$ for any state s and action a , $Q^{\pi_\alpha^*}(s, a) - c(s) > 0$. For such action a , we can establish a first order approximation of π_α^* in (18) based on

$$\begin{aligned} \pi_\alpha^*(s, a) \approx & 1 + \frac{1}{q-1} \left(\left(\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - \frac{c(s)}{\alpha} \right) \frac{q-1}{q} - 1 \right) \\ & + o \left(\left(\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - \frac{c(s)}{\alpha} \right) \frac{q-1}{q} - 1 \right) \end{aligned} \quad (21)$$

Now apply the constraint in (22) over all actions belonging to $S(s)$,

$$\sum_{a \in S(s)} \pi_\alpha^*(s, a) = 1 \quad (22)$$

we can obtain the result

$$c(s) \approx \alpha \frac{\sum_{a \in S(s)} \frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - q}{\|S(s)\|} + \alpha \left(q - \frac{q}{q-1} \right) \quad (23)$$

Clearly, when $q = 2$, $c(s)$ as approximated in (23) is identical to $c(s)$ in (19). Meanwhile, we can check the condition of $\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} > \frac{c(s)}{\alpha}$ whenever $a \in S(s)$. Apparently, only actions associated with high Q-values in state s have the chance to be performed by an RL agent. Suppose that $\{a_1, \dots, a_m\}$ are the actions with the m highest Q-values. For $S(s)$ to contain all these actions, we must make sure that

$$\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} > \frac{\sum_{i=1}^m \frac{Q^{\pi_\alpha^*}(s, a_i)}{\alpha} - q}{m} + \left(q - \frac{q}{q-1} \right) \quad (24)$$

Therefore,

$$m \frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} + q > \sum_{i=1}^m \frac{Q^{\pi_\alpha^*}(s, a_i)}{\alpha} + m \left(q - \frac{q}{q-1} \right) \quad (25)$$

Based on (25), $S(s)$ can now be estimated immediately as in (26).

$$S(s) \approx \left\{ a_i \mid \begin{array}{l} q + i \frac{Q^{\pi_\alpha^*}(s, a_i)}{\alpha} > \\ \sum_{j=1}^i \frac{Q^{\pi_\alpha^*}(s, a_j)}{\alpha} + i \left(q - \frac{q}{q-1} \right) \end{array} \right\} \quad (26)$$

Due to the inherent error involved in approximating $c(s)$ and $S(s)$ through (23) and (26) respectively, we do not know for sure which action can be safely ignored for future exploration. Without missing any potentially valuable actions, the second technique to approximate π_α^* is to use the *softplus function* [DBB⁺01] as a smooth implementation of $\max(\cdot, 0)$ in (18). Specifically,

$$\pi_\alpha^*(s, a) \propto \frac{1}{q-1} \delta \left(\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - \frac{c(s)}{\alpha} \right) \quad (27)$$

where the softplus function δ is defined as $\delta(x) = \log(1 + \exp(x))$. Based on (27), we can learn a Q-network parameterized by θ with the aim of minimizing the loss function $L^{\pi_\alpha^*}(\theta)$. $L^{\pi_\alpha^*}(\theta)$ is calculated based on (7) where the Bellman operator \mathcal{T}^* is replaced by $\mathcal{T}^{\pi_\alpha^*}$. Driven by this idea, we have developed an algorithm for Q-Learning under Tsallis entropy regularization as summarized in Algorithm 1.

3.4 Bellman Residue of Entropy Regularized Q-Learning

This subsection analyzes the Bellman residue to show that our Q-Learning algorithm will not suffer from any performance degradation despite of using approximated π_α^* in (27). We will particularly consider one extreme circumstance when $\alpha \rightarrow 0$. With α approaching to 0, it is straightforward to verify that only the action that produces the highest Q-value in any state s , i.e. a_1 , satisfies the condition in (26). Therefore $c(s) = \alpha \left(\frac{Q^{\pi_\alpha^*}(s, a_1)}{\alpha} - \frac{q}{q-1} \right)$. According to (27), $\pi_\alpha^*(s, a_1) \propto \delta(-\frac{q}{q-1}) > 0$. On the other hand, for any action $a \neq a_1$,

$$\begin{aligned} & \lim_{\alpha \rightarrow 0} {}^{q-1}\sqrt{\delta \left(\frac{Q^{\pi_\alpha^*}(s, a)}{\alpha} - \frac{c(s)}{\alpha} \right)} \\ &= \lim_{\alpha \rightarrow 0} {}^{q-1}\sqrt{\delta \left(\frac{Q^{\pi_\alpha^*}(s, a) - Q^{\pi_\alpha^*}(s, a_1)}{\alpha} + \frac{q}{q-1} \right)} \\ &= 0 \end{aligned} \tag{28}$$

We hence can conclude that, when $\alpha \rightarrow 0$, $\pi_\alpha^*(s, a_1) \rightarrow 1$. Consequently, π_α^* degenerates to the value-maximization policy. Based on this understanding, we can further prove Theorem 1 below (see Appendix for proof).

Theorem 1 *For the Q-learning problem under Tsallis entropy regularization, suppose that π_α^* is the approximated optimal policy for the problem as defined in (27) and α is the corresponding regularization coefficient, then for any state $s \in \mathbb{S}$ and any action $a \in \mathbb{A}$, the Bellman residue $|\mathcal{T}^* Q^{\pi_\alpha^*} - Q^{\pi_\alpha^*}|$ satisfies the following property*

$$\lim_{\alpha \rightarrow 0} \left| \mathcal{T}^* Q^{\pi_\alpha^*}(s, a) - Q^{\pi_\alpha^*}(s, a) \right| = 0$$

Because the Bellman residue converges to 0 with decreasing α , when α is sufficiently small, the performance of Algorithm 1 is as good as DQN and Double DQN. For this reason, Algorithm 1 employs a linear schedule to constantly decrease α , thereby gradually reducing entropy-induced exploration till $\alpha = 0$ and learning converges.

4 Bootstrapped Q-Ensemble under Tsallis Entropy Regularization

Using Algorithm 1 alone is insufficient to realize deep and effective exploration. Following the *deep exploration principle* studied in [ORW14, ORWR17], we expand an MDP \mathcal{M}_q in this paper with the Tsallis entropy regularizer in (13) under the settings of $1 < q < q_{max} < \infty$. In other words, the immediate reward of performing any action in state s by following policy π is extended with a new term that depends on Tsallis entropy of π in state s . Given the experiences obtained so far by an RL agent through direct interactions with its learning environment, denoted as \mathcal{B} , a posterior model over \mathcal{M}_q can be established in theory and represented as $P(\mathcal{M}_q|\mathcal{B})$.

Deep exploration requires an RL agent to randomly sample one MDP \mathcal{M}_q from $P(\mathcal{M}_q|\mathcal{B})$ and subsequently utilize $Q^{\pi_{\alpha, \mathcal{M}_q}^*}$ and $\pi_{\alpha, \mathcal{M}_q}^*$ to control future interactions with its learning environment in the next problem episode. Each problem episode starts from an initial state s_0 and ends whenever a final state is reached. $Q^{\pi_{\alpha, \mathcal{M}_q}^*}$ and $\pi_{\alpha, \mathcal{M}_q}^*$ stand respectively for the optimal

Q-function and optimal policy with respect to an MDP \mathcal{M}_q under the specific settings of q and α . Apparently this deep exploration method has the aim of optimizing the posterior learning performance of an RL agent based on its past experiences, as described below

$$J^*(\mathcal{B}) = \mathbb{E}_{\mathcal{M}_q \sim P(\mathcal{M}_q|\mathcal{B})} \mathbb{E}_{(s,a) \sim \pi_{\alpha, \mathcal{M}_q}^*} \left(r(s,a) + \alpha H_q^{\pi_{\alpha, \mathcal{M}_q}^*}(s) \right) \quad (29)$$

On large-scale RL problems it is difficult to keep track of $P(\mathcal{M}_q|\mathcal{B})$ as well as to determine the optimal policy for every possible \mathcal{M}_q . Inspired by [OBPR16], we decide to efficiently approximate the deep exploration process through bootstrapping. This is implemented in the BQETR algorithm (see Algorithm 2) by maintaining an ensemble of entropy regularized Q-networks. All Q-networks share the same regularization coefficient α . Meanwhile, different Q-networks follow different settings of q so as to enhance the diversity of the ensemble, which is essential for effective deep exploration.

Apparently, no change to q is required for any Q-network in the ensemble during RL. Since agent’s past experiences \mathcal{B} will not affect the posterior distribution over q , each Q-network in the ensemble can be sampled equally likely for the next episode of deep exploration. This simple technique enables us to seamlessly integrate entropy-induced exploration with bootstrap-induced exploration and lays the foundation of the BQETR algorithm. Because α is decremented each time by a very small step Δ_α in Algorithm 1 and Algorithm 2, its change will not affect the effectiveness of the bootstrapping mechanism.

Algorithm 2 The Bootstrapped Q-Ensemble under Tsallis Entropy Regularization (BQETR) Algorithm

- 1: **Input:** an ensemble of K Q-networks $\{Q_k\}_{k=1}^K$, a list of q values $\{q_k\}_{k=1}^K$ for the Tsallis entropy regularizers, α_0 for the initial regularization coefficient, a replay buffer \mathcal{B} that stores past state-transition samples for training, and a masking distribution M .
 - 2: **for** each problem episode **do**:
 - 3: Choose the i -th Q-network in $\{Q_k\}_{k=1}^K$ randomly
 - 4: Obtain initial state s_0 from environment
 - 5: **for** $t = 1, \dots$ until end of episode **do**:
 - 6: Use Q_i, q_i and α to sample action a_t according to (27)
 - 7: Perform a_t
 - 8: Sample bootstrap mask $m_t \sim M$
 - 9: Add $(s_t, a_t, s_{t+1}, r_t, m_t)$ to \mathcal{B}
 - 10: **if** learning interval is reached **do**:
 - 11: Follow Algorithm 1 to train all K Q-networks.
 - 12: Reduce α linearly by Δ_α until 0
-

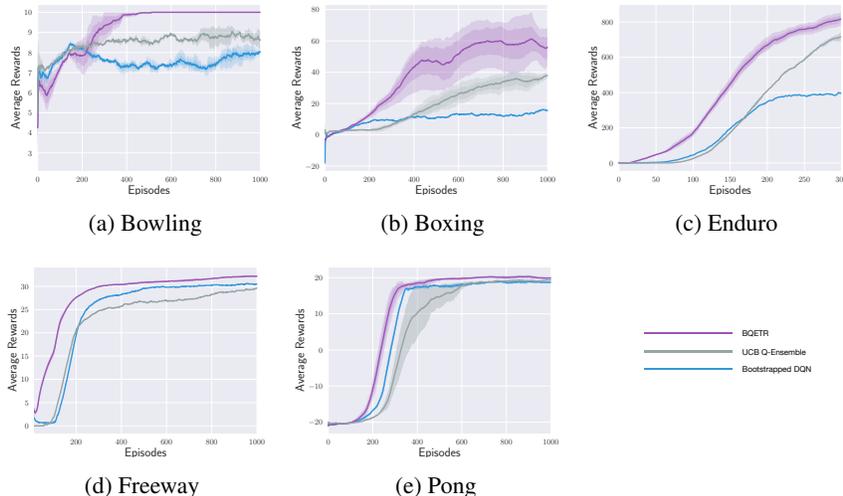


Figure 1: Average total return per episode obtained by BQETR, Bootstrapped DQN and UCB Q-Ensemble on five Atari game playing tasks, including Bowling, Boxing, Enduro, Freeway, and Pong.

5 Experiment

In this section, the learning performance of BQETR is compared to the performance achievable through Bootstrapped DQN and UCB Q-Ensemble on commonly studied Atari game playing tasks [HGS16, MKS⁺15, SLA⁺15, OBPR16, CSAS17]. Based on the performance results, the sample complexity of the three algorithms is further analyzed to demonstrate that BQETR can improve sample efficiency through deep and effective exploration of its learning environment.

In this paper we consider specifically five video games simulated by the Arcade Learning Environment [BNVB15] as benchmark problems, including Bowling, Boxing, Enduro, Freeway and Pong. These problems require an RL agent to handle high-dimensional state spaces (i.e. an agent must be able to process direct video input provided by the games) and are highly difficult to solve, even for expert human game players. As a result, to achieve reasonable learning performance, an RL agent must play numerous rounds of each benchmark game. Hence they are suitable problems to reveal the difference in sample efficiency upon using various exploration methods for DRL.

We implement all algorithms in the experiments based on the high-quality implementation of Double DQN provided by OpenAI baselines [DHK⁺17]. We also closely follow the parametric settings of Bootstrapped DQN and UCB Q-Ensemble presented in [OBPR16, CSAS17]. Meanwhile, for a fair comparison, we use identical settings for common parameters shared between BQETR and Bootstrapped DQN. BQETR also introduces two additional parameters, i.e. the initial value for the entropy regularization coefficient α_0 and the *entropic index* q for Tsallis entropy. Without spending substantial efforts in fine-tuning these parameters, α_0 is set to 0.5 (other settings ranging from 1.0 to 0.1 do not seem to produce noticeable difference in performance). After each learning interval, α will be decremented by 0.5×10^{-5} till 0. Since the Q-ensemble maintained by all algorithms contains 10 individual Q-networks. The values for q in BQETR have been set to 1.5, 1.6, \dots , 2.4 respectively for each Q-network. Moreover, every

pixel input to a Q-network is obtained by averaging the same pixel over four consecutive frames of the game video. On each game playing task, we have run every algorithm for only 3M frames by using commodity desktop computers (no GPUs). This enables us to examine the effectiveness of all algorithms under limited computation resources and sample budget.

5.1 Results on Learning Effectiveness

Figure 1 depicts the learning performance (i.e., average total return per episode) of the three algorithms in our experiments. To cover at least 3M frames, the performance across 1000 learning episodes have been presented in the figure, except for Enduro. This is because each episode in Enduro includes more frames. Therefore 3M frames have been reached after just playing 300 episodes of Enduro.

As evidenced in Figure 1, BQETR outperformed Bootstrapped DQN and UCB Q-Ensemble on all benchmark problems. Particularly on Bowling, Boxing and Enduro, BQETR achieved significantly higher performance than competing algorithms. In the meantime, BQETR also managed to solve Freeway and Pong clearly faster than other algorithms. Based on the experiment results, we believe that BQETR is an effective algorithm for DRL thanks to its integrated use of both entropy-induced and bootstrap-induced exploration techniques.

5.2 Results on Sample Efficiency

To analyze sample efficiency, we adopt the performance metrics introduced in [SWD⁺17]. Particularly, the fast learning metric in Table 1 calculates the average performance across all learning episodes and the final performance metric calculates the average performance obtained in the last 10 episodes. Both metrics in Table 1 clearly show that BQETR is not only effective in terms of final performance but also significantly more sample efficient than both Bootstrapped DQN and UCB Q-Ensemble.

Scoring Metric	Algorithms	Bowling	Boxing	Enduro	Freeway	Pong
Fast Learning	BQETR	9.36 ± 0.47	60.61 ± 17.71	486.13 ± 34.56	29.33 ± 0.63	12.76 ± 0.41
	Bootstrapped DQN	7.63 ± 0.55	19.80 ± 1.00	279.98 ± 10.00	26.19 ± 0.30	11.47 ± 0.24
	UCB Q-Ensemble	8.59 ± 0.41	38.11 ± 6.36	390.68 ± 24.10	25.08 ± 0.33	10.72 ± 1.32
Final Performance	BQETR	10.00 ± 0.02	57.41 ± 7.11	812.10 ± 36.39	32.22 ± 0.28	20.74 ± 0.38
	Bootstrapped DQN	7.97 ± 0.73	41.50 ± 1.00	374.23 ± 5.00	30.86 ± 0.10	19.23 ± 0.18
	UCB Q-Ensemble	9.20 ± 0.09	47.96 ± 2.34	725.95 ± 30.50	31.00 ± 0.10	19.34 ± 0.30

Table 1: Scoring metrics of fast learning and final performance obtained by BQETR, Bootstrapped DQN and UCB Q-Ensemble on five Atari games.

6 Conclusions

In this paper we studied entropy-induced environment exploration via deep Q-learning under general Tsallis entropy regularization. Through this study, we developed the first time in literature new approximation techniques to address entropy regularized RL problems. Bellman residue analysis subsequently showed that our approximation techniques will not affect the final performance achievable through Q-learning. Driven by the goal for deep exploration, we have further developed a bootstrapped Q-learning algorithm involving an ensemble of Q-networks. Every Q-network is controlled by a Tsallis entropy regularizer under different settings of q so as to achieve high ensemble diversity and effective deep exploration.

Looking into the future, it is interesting to explore the possibilities of extending our Q-learning algorithm to tackle RL problems with high-dimensional and continuous action spaces. Meanwhile, it is also interesting to explore the benefits of our Q-learning algorithm on more problem domains. Due to limited computation resources that are available to this research, we cannot conduct large-scale experimental studies in this paper. However our experiment results have clearly shown that our new algorithm is both effective and sample efficient.

Appendix

This appendix presents proof of Theorem 1, i.e. $\|\mathcal{T}^*Q^{\pi_\alpha^*} - Q^{\pi_\alpha^*}\| \rightarrow 0$ by decreasing the regularization coefficient α all the way to 0. Specifically, for any state-action pair (s, a) , we can derive the following inequalities.

$$\begin{aligned}
0 &\leq \left| \mathcal{T}^*Q^{\pi_\alpha^*}(s, a) - \mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) \right| \\
&\leq \mathbb{E}_{s' \sim P(s, a, s')} \left| \max_b Q^{\pi_\alpha^*}(s', b) - \sum_b \pi_\alpha^*(s', b) Q^{\pi_\alpha^*}(s', b) \right| \\
&\leq \mathbb{E}_{s' \sim P(s, a, s')} \left(\sum_b |\mathbb{I}_{b=b_1}^{s'} - \pi_\alpha^*(s', b)| \cdot |Q^{\pi_\alpha^*}(s', b)| \right) \\
&\leq \mathbb{E}_{s' \sim P(s, a, s')} \left(\sum_b |\mathbb{I}_{b=b_1}^{s'} - \pi_\alpha^*(s', b)| \sum_b |Q^{\pi_\alpha^*}(s', b)| \right)
\end{aligned} \tag{30}$$

where $\mathbb{I}_{b=b_1}^{s'}$ refers to the policy that selects action b_1 in state s' with probability 1 and b_1 is the action with the highest Q-value in state s' . Assume without loss of generality that the absolute Q-value with respect to any state and any action can never exceed \bar{Q} . Then, from (30), we have

$$\begin{aligned}
&\left| \mathcal{T}^*Q^{\pi_\alpha^*}(s, a) - \mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) \right| \\
&\leq \|\mathbb{A}\| \cdot \bar{Q} \cdot \mathbb{E}_{s' \sim P(s, a, s')} \left(2D_{TV} \left(\mathbb{I}_{b=b_1}^{s'} \|\pi_\alpha^*(s', \cdot)\| \right) \right)
\end{aligned} \tag{31}$$

with $D_{TV}(x||y) = \frac{1}{2} \sum_i |x_i - y_i|$ representing the *total variation divergence* in between any two discrete probability distributions x and y . Due to the fact that $D_{TV}(x||y)^2 \leq D_{KL}(x||y)$ where D_{KL} is the standard *KL divergence* [Pol], we can further obtain the inequality below from (31),

$$\begin{aligned}
&\left| \mathcal{T}^*Q^{\pi_\alpha^*}(s, a) - \mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) \right| \\
&\leq 2\|\mathbb{A}\| \cdot \bar{Q} \cdot \mathbb{E}_{s' \sim P(s, a, s')} \sqrt{D_{KL} \left(\mathbb{I}_{b=b_1}^{s'} \|\pi_\alpha^*(s', \cdot)\| \right)} \\
&= 2\|\mathbb{A}\| \cdot \bar{Q} \cdot \mathbb{E}_{s' \sim P(s, a, s')} \sqrt{-\log \pi_\alpha^*(s', b_1)}
\end{aligned} \tag{32}$$

Consequently,

$$\begin{aligned}
0 &\leq \lim_{\alpha \rightarrow 0} \left| \mathcal{T}^*Q^{\pi_\alpha^*}(s, a) - \mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) \right| \\
&\leq 2\|\mathbb{A}\| \cdot \bar{Q} \cdot \mathbb{E}_{s' \sim P(s, a, s')} \lim_{\alpha \rightarrow 0} \sqrt{-\log \pi_\alpha^*(s', b_1)} \\
&= 0
\end{aligned} \tag{33}$$

Using (33), it can be further shown that

$$\begin{aligned}
0 &\leq \lim_{\alpha \rightarrow 0} \left| \mathcal{T}^* Q^{\pi_\alpha^*}(s, a) - Q^{\pi_\alpha^*}(s, a) \right| \\
&\leq \lim_{\alpha \rightarrow 0} \left| \mathcal{T}^* Q^{\pi_\alpha^*}(s, a) - \mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) \right| \\
&\quad + \lim_{\alpha \rightarrow 0} \left| \mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) - Q^{\pi_\alpha^*}(s, a) \right| \\
&= 0
\end{aligned} \tag{34}$$

Notice that when $\alpha \rightarrow 0$, the error involved in approximating π_α^* in (18) is negligible since the action that produces the highest Q-value will be selected with probability 1. As a result, $\mathcal{T}^\alpha Q^{\pi_\alpha^*}(s, a) - Q^{\pi_\alpha^*}(s, a) = 0$ for any state s and action a .

References

- [AO07] P. Auer and R. Ortner. Logarithmic online regret bounds for undiscounted reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 49–56, 2007.
- [Ber95] D. P. Bertsekas. *Dynamic programming and optimal control*. Athena scientific Belmont, MA, 1995.
- [BNVB15] Marc G. Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. In *Proceedings of the 24th International Conference on Artificial Intelligence, IJCAI’15*, pages 4148–4152. AAAI Press, 2015.
- [BT02] R. I. Brafman and M. Tennenholtz. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct):213–231, 2002.
- [CPZ18a] G. Chen, Y. Peng, and M. Zhang. An adaptive clipping approach for proximal policy optimization. *arXiv preprint arXiv:1804.06461*, 2018.
- [CPZ18b] G. Chen, Y. Peng, and M. Zhang. Constrained expectation-maximization methods for effective reinforcement learning. In *International Joint Conference on Neural Networks*, 2018.
- [CSAS17] R. Y. Chen, S. Sidor, P. Abbeel, and J. Schulman. UCB Exploration via Q-Ensembles. *arXiv preprint arXiv:1706.01502*, 2017.
- [DB15] C. Dann and E. Brunskill. Sample complexity of episodic fixed-horizon reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 2818–2826, 2015.
- [DBB⁺01] C. Dugas, Y. Bengio, F. Bélisle, C. Nadeau, and R. Garcia. Incorporating second-order functional knowledge for better option pricing. In *Advances in neural information processing systems*, pages 472–478, 2001.
- [DHK⁺17] Prafulla Dhariwal, Christopher Hesse, Oleg Klimov, Alex Nichol, Matthias Plappert, Alec Radford, John Schulman, Szymon Sidor, and Yuhuai Wu. Openai baselines. <https://github.com/openai/baselines>, 2017.
- [GMP⁺15] M. Ghavamzadeh, S. Mannor, J. Pineau, A. Tamar, et al. Bayesian reinforcement learning: A survey. *Foundations and Trends in Machine Learning*, 8(5-6):359–483, 2015.

- [HGS16] H. Van Hasselt, A. Guez, and D. Silver. Deep Reinforcement Learning with Double Q-Learning. In *Thirtieth AAAI Conference on Artificial Intelligence*, volume 16, pages 2094–2100, 2016.
- [HTAL17] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine. Reinforcement learning with deep energy-based policies. *arXiv preprint arXiv:1702.08165*, 2017.
- [Kak03] S. Kakade. *On the sample complexity of reinforcement learning*. PhD thesis, University College London, 2003.
- [KS02] M. Kearns and S. Singh. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49(2-3):209–232, 2002.
- [LCO17] K. Lee, S. Choi, and S. Oh. Sparse Markov Decision Processes with Causal Sparse Tsallis Entropy Regularization for Reinforcement Learning. *arXiv preprint arXiv:1709.06293*, 2017.
- [LHP⁺15] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [LLW16] L. Li, Y. Lv, and F. Y. Wang. Traffic signal timing via deep reinforcement learning. *IEEE/CAA Journal of Automatica Sinica*, 3(3):247–254, 2016.
- [MKS⁺15] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland and G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [NCG18] O. Nachum, Y. Chow, and M. Ghavamzadeh. Path Consistency Learning in Tsallis Entropy Regularized MDPs. *arXiv preprint arXiv:1802.03501*, 2018.
- [NNXS17] O. Nachum, M. Norouzi, K. Xu, and D. Schuurmans. Bridging the gap between value and policy based reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 2772–2782, 2017.
- [OBPR16] I. Osband, C. Blundell, A. Pritzel, and B. Van Roy. Deep exploration via bootstrapped DQN. In *Advances in neural information processing systems*, pages 4026–4034, 2016.
- [OMKM17] B. O’Donoghue, R. Munos, K. Kavukcuoglu, and V. Mnih. Combining policy gradient and Q-learning. *arXiv preprint arXiv:1611.01626*, 2017.
- [OR16] I. Osband and B. Van Roy. Why is posterior sampling better than optimism for reinforcement learning. *arXiv preprint arXiv:1607.00215*, 2016.
- [ORR13] I. Osband, D. Russo, and B. Van Roy. (More) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems*, pages 3003–3011, 2013.
- [ORW14] I. Osband, B. Van Roy, and Z. Wen. Generalization and exploration via randomized value functions. *arXiv preprint arXiv:1402.0635*, 2014.
- [ORWR17] I. Osband, D. Russo, Z. Wen, and B. Van Roy. Deep exploration via randomized value functions. *arXiv preprint arXiv:1703.07608*, 2017.
- [Pol] D. Pollard. *Asymptopia: an exposition of statistical asymptotic theory*. 2000. URL <http://www.stat.yale.edu/~pollard/Books/Asymptopia>.
- [PP93] A. R. Plastino and A. Plastino. Tsallis’ entropy, ehrenfest theorem and information theory. *Physics Letters A*, 177(3):177–179, 1993.

- [RR13] D. Russo and B. Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems*, pages 2256–2264, 2013.
- [RR14] D. Russo and B. Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- [SAC17] J. Schulman, P. Abbeel, and X. Chen. Equivalence between policy gradients and soft Q-Learning. *arXiv preprint arXiv:1704.06440*, 2017.
- [SLA⁺15] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, pages 1889–1897, 2015.
- [SLW⁺06] A. L. Strehl, L. Li, E. Wiewiora, J. Langford, and M. L. Littman. Pac model-free reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 881–888. ACM, 2006.
- [SMSM00] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063, 2000.
- [Str07] A. L. Strehl. *Probably approximately correct (PAC) exploration in reinforcement learning*. PhD thesis, Rutgers University-Graduate School-New Brunswick, 2007.
- [SWD⁺17] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *CoRR*, 2017.
- [Tsa94] C. Tsallis. Nonextensive physics: a possible connection between generalized statistical mechanics and quantum groups. *Physics Letters A*, 195(5-6):329–334, 1994.
- [WBH⁺16] Z. Wang, V. Bapst, N. Heess, V. Mnih, R. Munos, K. Kavukcuoglu, and N. de Freitas. Algorithms for multi-armed bandit problems. *arXiv preprint arXiv:1611.01224*, 2016.
- [WMG⁺17] Y. Wu, E. Mansimov, R. B. Grosse, S. Liao, and J. Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *Advances in neural information processing systems*, pages 5279–5288, 2017.