

---

# Implicitly Regularized RL with Implicit Q-values

---

Nino Vieillard<sup>1,2</sup>

Marcin Andrychowicz<sup>1</sup>

Anton Raichuk<sup>1</sup>

Olivier Pietquin<sup>1</sup>

Matthieu Geist<sup>1</sup>

<sup>1</sup>Google Research, Brain Team    <sup>2</sup>Université de Lorraine, CNRS, Inria, IECL, F-54000 Nancy, France

## Abstract

The  $Q$ -function is a central quantity in many Reinforcement Learning (RL) algorithms for which RL agents behave following a (soft)-greedy policy w.r.t. to  $Q$ . It is a powerful tool that allows action selection without a model of the environment and even without explicitly modeling the policy. Yet, this scheme can only be used in discrete action tasks, with small numbers of actions, as the softmax over actions cannot be computed exactly otherwise. More specifically, the usage of function approximation to deal with continuous action spaces in modern actor-critic architectures intrinsically prevents the exact computation of a softmax. We propose to alleviate this issue by parametrizing the  $Q$ -function *implicitly*, as the sum of a log-policy and a value function. We use the resulting parametrization to derive a practical off-policy deep RL algorithm, suitable for large action spaces, and that enforces the softmax relation between the policy and the  $Q$ -value. We provide a theoretical analysis of our algorithm: from an Approximate Dynamic Programming perspective, we show its equivalence to a regularized version of value iteration, accounting for both entropy and Kullback-Leibler regularization, and that enjoys beneficial error propagation results. We then evaluate our algorithm on classic control tasks, where its results compete with state-of-the-art methods.

## 1 INTRODUCTION

A large body of reinforcement learning (RL) algorithms, based on approximate dynamic programming

(ADP) (Bertsekas and Tsitsiklis, 1996; Scherrer et al., 2015), operate in two steps: a greedy step, where the algorithm learns a policy that maximizes a  $Q$ -value, and an evaluation step, that (partially) updates the  $Q$ -values towards the  $Q$ -values of the policy. A common improvement to these techniques is to use regularization, that prevents the new updated policy from being too different from the previous one, or from a fixed “prior” policy. For example, Kullback-Leibler (KL) regularization keeps the policy close to the previous iterate (Vieillard et al., 2020a), while entropy regularization keeps the policy close to the uniform one (Haarnoja et al., 2018a). Entropy regularization, often used in this context (Ziebart, 2010), modifies both the greedy step and the evaluation step so that the policy jointly maximizes its expected return and its entropy. In this framework, the solution to the policy optimization step is simply a softmax of the  $Q$ -values over the actions. In small discrete action spaces, the softmax can be computed exactly: one only needs to define a *critic* algorithm, with a single loss that optimizes a  $Q$ -value. However, in large multi-dimensional – or even continuous – action spaces, one needs to estimate it. This estimation is usually done by adding an *actor* loss, that optimizes a policy to fit this softmax. It results in an (off-policy) *actor-critic* algorithm, with two losses that are optimized simultaneously<sup>1</sup> (Degris et al., 2012). This additional optimization step introduces supplementary errors to the ones already created by the approximation in the evaluation step.

To remove these extraneous approximations, we introduce the Implicit  $Q$ -values (IQ) algorithm, that deviates from classic actor-critics, as it optimizes a policy and a value in a single loss. The core idea is to implicitly represent the  $Q$ -value as the sum of a value function and a log-policy. This representation ensures that the policy is the *exact* softmax of the  $Q$ -value, *despite the use of any approximation scheme*. We use this to design a practical model-free deep RL algorithm

---

Proceedings of the 25<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2022, Valencia, Spain. PMLR: Volume 151. Copyright 2022 by the author(s).

---

<sup>1</sup>We refer here specifically to *off-policy* actor-critics, built on a value-iteration-like scheme. For on-policy actor-critics, losses are optimized sequentially, and the policy is usually fully evaluated.

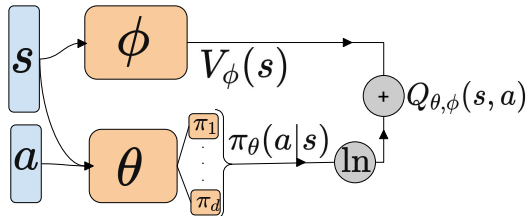


Figure 1: view of the IQ parametrization.

that optimizes with a single loss a policy network and a value network, built on this implicit representation of a  $Q$ -value. To better understand it, we abstract this algorithm to an ADP scheme, IQ-DP, and use this point of view to provide a detailed theoretical analysis. It relies on a key observation, that shows an equivalence between IQ-DP and a specific form of regularized Value Iteration (VI). This equivalence explains the role of the components of IQ: namely, IQ performs entropy and KL regularization. It also allows us to derive strong performance bounds for IQ-DP. In particular, we show that the errors made when following IQ-DP are compensated along iterations.

Parametrizing the  $Q$ -value as a sum of a log-policy and a value is reminiscent of the dueling architecture (Wang et al., 2016), that factorizes the  $Q$ -value as the sum of an advantage and a value. In fact, we show that it is a limiting case of IQ in a discrete actions setting. This link highlights the role of our policy, which calls for a discussion on its necessary parametrization.

Finally, we empirically validate IQ. We evaluate our method on several classic continuous control benchmarks: locomotion tasks from Openai Gym (Brockman et al., 2016), and hand manipulation tasks from the Adroit environment (Rajeswaran et al., 2017). On these environments, IQ reaches performances competitive with state-of-the-art actor-critic methods.

## 2 IMPLICIT $Q$ -VALUE PARAMETRIZATION

We consider the standard Reinforcement Learning (RL) setting, formalized as a Markov Decision Process (MDP). An MDP is a tuple  $\{\mathcal{S}, \mathcal{A}, P, r, \gamma\}$ .  $\mathcal{S}$  and  $\mathcal{A}$  are the finite state and action spaces<sup>2</sup>,  $\gamma \in [0, 1)$  is the discount factor and  $r : \mathcal{S} \times \mathcal{A} \rightarrow [-R_{max}, R_{max}]$  is the bounded reward function. Write  $\Delta_X$  the simplex over the finite set  $X$ . The dynamics of an MDP are defined by a Markovian transition kernel  $P \in \Delta_{\mathcal{S} \times \mathcal{A}}$ , where  $P(s'|s, a)$  is the probability of transitioning to state  $s'$  after taking action  $a$  in  $s$ . An RL agent acts

<sup>2</sup>We restrict to finite spaces for the sake of analysis, but our approach applies to continuous spaces.

through a stationary stochastic policy  $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ , a mapping from states to distribution over actions. The quality of a policy is quantified by the value function,  $V_{\pi}(s) = \mathbb{E}_{\pi}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) | s_0 = s]$ . The  $Q$ -function is a useful extension, which notably allows choosing a (soft)-greedy action in a model-free setting,  $Q_{\pi}(s, a) = r(s, a) + \mathbb{E}_{s'|s, a}[V_{\pi}(s')]$ . An optimal policy is one that achieve the highest expected return,  $\pi_* = \operatorname{argmax}_{\pi} V_{\pi}$ .

A classic way to design practical algorithms beyond the tabular setting is to adopt the Actor-Critic perspective. In this framework, an RL agent parametrizes a policy  $\pi_{\theta}$  and a  $Q$ -value  $Q_{\psi}$  with function approximation, usually through the use of neural networks, and aims at estimating an optimal policy. The policy and the  $Q$ -function are then updated by minimizing two losses: the actor loss corresponds to the greedy step, and the critic loss to the evaluation step. The weights of the policy and  $Q$ -value networks are regularly frozen into *target* weights  $\bar{\psi}$  and  $\bar{\theta}$ . With entropy regularization, the greedy step amounts to finding the policy that maximizes  $\mathbb{E}_{s \sim \mathcal{S}, a \sim \pi_{\theta}}[Q_{\bar{\psi}}(s, a) - \tau \ln \pi_{\theta}(a|s)]$  (maximize the  $Q$ -value with stochastic enough policy). The solution to this problem is simply  $\pi_{\theta}(\cdot|s) = \operatorname{softmax}(Q_{\bar{\psi}}(s, \cdot)/\tau)$ , which is the result of the greedy step of regularized Value Iteration (VI) (Geist et al., 2019) and, for example, how the optimization step of Soft Actor-Critic (Haarnoja et al., 2018a, SAC) is built. In a setting where the action space is discrete and small, it amounts to a simple softmax computation. However, on more complex action spaces (continuous, and/or with a higher number of dimensions: as a reference, the Humanoid-v2 environment from Openai Gym (Brockman et al., 2016) has an action space of dimension 17), it becomes prohibitive to use the exact solution. In this case, the common practice is to resort to an approximation with a parametric distribution model. In many actor critic algorithms (SAC, MPO (Abdolmaleki et al., 2018), ...), the policy is modelled as a Gaussian distribution over actions. It introduces approximation errors, resulting from the partial optimization process of the critic, and inductive bias, as a Gaussian policy cannot represent an arbitrary softmax distribution. We now turn to the description of our core contribution: the Implicit  $Q$ -value (IQ) algorithm, introduced to mitigate this discrepancy.

IQ implicitly parametrizes a  $Q$ -value via an explicit parametrization of a policy and a value, as visualized in Fig. 1. Precisely, from a policy network  $\pi_{\theta}$  and a value network  $V_{\phi}$ , we define our implicit  $Q$ -value as

$$Q_{\theta, \phi}(s, a) = \tau \ln \pi_{\theta}(a|s) + V_{\phi}(s). \quad (1)$$

Since  $\pi_{\theta}$  is constrained to be a distribution over the actions, we have by construction that  $\pi_{\theta}(a|s) =$

$\text{softmax}(Q_{\theta,\phi}/\tau)$ , the solution of the regularized greedy step (see Appx. A.1 for a detailed proof). Hence, the consequence of using such a parametrization is that the greedy step is performed exactly, even in the function approximation regime. Compared to the classic actor-critic setting, it thus gets rid of the errors created by the actor. Note that calling  $V_\phi$  a value makes sense, since following the same reasoning we have that  $V_\phi(s) = \tau \ln \sum_{a'} \exp(Q_{\theta,\phi}(s, a')/\tau)$ , a soft version of the value. With this parametrization in mind, one could derive a deep RL algorithm from any value-based loss using entropy regularization. We conserve the fixed-point approach of the standard actor-critic framework,  $\theta$  and  $\phi$  are regularly copied to  $\bar{\theta}$  and  $\bar{\phi}$ , and we design an off-policy algorithm, working on a replay buffer of transitions  $(s_t, a_t, r_t, s_{t+1})$  collected during training. Consider two hyperparameters,  $\tau \in (0, \infty)$  and  $\alpha \in (0, 1)$  that we will show in Sec. 3 control two forms of regularization. The policy and value are optimized jointly by minimizing the loss

$$\mathcal{L}_{\text{IQ}}(\theta, \phi) = \hat{\mathbb{E}} \left[ \left( r_t + \alpha \tau \ln \pi_{\bar{\theta}}(a_t | s_t) + \gamma V_{\bar{\phi}}(s_{t+1}) - \tau \ln \pi_{\theta}(a_t | s_t) - V_{\phi}(s_t) \right)^2 \right], \quad (2)$$

where  $\hat{\mathbb{E}}$  denote the empirical expected value over a dataset of transitions. IQ consists then in a single loss that optimizes jointly a policy and a value. This brings a notable remark on the role of  $Q$ -functions in RL. Indeed,  $Q$ -learning was introduced by [Watkins and Dayan \(1992\)](#) – among other reasons – to make greediness possible without a model (using a value only, one needs to maximize it over all possible successive states, which requires knowing the transition model), and consequently derive practical, model-free RL algorithms. Here however, IQ illustrates how, with the help of regularization, one can derive a model-free algorithm that does not rely on an explicit  $Q$ -value.

### 3 ANALYSIS

In this section, we explain the workings of the IQ algorithm defined by Eq. (2) and detail the influence of its hyperparameters. We abstract IQ into an ADP framework, and show that, from that perspective, it is equivalent to a Mirror Descent VI (MD-VI) scheme ([Geist et al., 2019](#)), with both entropy and KL regularization. Let us first introduce some useful notations. We make use of the actions partial dot-product notation: for  $u, v \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ , we define  $\langle u, v \rangle = \left( \sum_{a \in \mathcal{A}} u(s, a)v(s, a) \right)_s \in \mathbb{R}^{\mathcal{S}}$ . For any  $V \in \mathbb{R}^{\mathcal{S}}$ , we have for any  $(s, a) \in \mathcal{S} \times \mathcal{A}$   $PV(s, a) = \sum_{s'} P(s' | s, a)V(s')$ . We will define regularized algorithms, using the entropy of a policy,  $\mathcal{H}(\pi) = -\langle \pi, \ln \pi \rangle$ , and the KL divergence between two policies,  $\text{KL}(\pi || \mu) = \langle \pi, \ln \pi - \ln \mu \rangle$ .

The  $Q$ -value of a policy is the unique fixed point of its Bellman operator  $T_\pi$  defined for any  $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  as  $T_\pi Q = r + \gamma P \langle \pi, Q \rangle$ . We denote  $Q_* = Q_{\pi_*}$  the optimal  $Q$ -value (the  $Q$ -value of the optimal policy). When the MDP is entropy-regularized with a temperature  $\tau$ , a policy  $\pi$  admits a *regularized*  $Q$ -value  $Q_\pi^\tau$ , the fixed point of the regularized Bellman operator  $T_\pi^\tau Q = r + \gamma P \langle \pi, Q - \tau \ln \pi \rangle$ . A regularized MDP admits an optimal *regularized* policy  $\pi_*^\tau$  and a unique optimal *regularized*  $Q$ -value  $Q_*^\tau$  ([Geist et al., 2019](#)).

#### 3.1 Ideal case

First, let us look at the ideal case, *i.e.* when  $\mathcal{L}_{\text{IQ}}$  is exactly minimized at each iteration (tabular representation, dataset covering the whole state-action space, expectation rather than sampling for transitions). In this context, IQ can be understood as a Dynamic Programming (DP) scheme that iterates on a policy  $\pi_{k+1}$  and a value  $V_k$ . They are respectively equivalent to the target networks  $\pi_{\bar{\theta}}$  and  $V_{\bar{\phi}}$ , while the next iterate  $(\pi_{k+2}, V_{k+1})$  matches the solution  $(\pi_\theta, V_\phi)$  of the optimization problem in Eq. (2). We call the scheme IQ-DP( $\alpha, \tau$ ) and one iteration is defined by choosing  $(\pi_{k+2}, V_{k+1})$  such that the squared term in Eq. (2) is 0,

$$\tau \ln \pi_{k+2} + V_{k+1} = r + \alpha \tau \ln \pi_{k+1} + \gamma PV_k. \quad (3)$$

This equation is well-defined, due to the underlying constraint that  $\pi_{k+2} \in \Delta_{\mathcal{A}}^{\mathcal{S}}$  (the policy must be a distribution over actions), that is  $\sum_{a \in \mathcal{A}} \pi(a|s) = 1$  for all  $s \in \mathcal{S}$ . The basis for our discussion will be the equivalence of this scheme to a version of regularized VI. Indeed, we have the following result, proved in Appendix A.3.

**Theorem 1.** *For any  $k \geq 1$ , let  $(\pi_{k+2}, V_{k+1})$  be the solution of IQ-DP( $\alpha, \tau$ ) at step  $k$ . We have that*

$$\begin{cases} \pi_{k+2} = \text{argmax} \langle \pi, r + \gamma PV_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi) \\ \quad \quad \quad - \alpha \tau \text{KL}(\pi || \pi_{k+1}) \\ V_{k+1} = \langle \pi_{k+2}, r + \gamma PV_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi_{k+2}) \\ \quad \quad \quad - \alpha \tau \text{KL}(\pi_{k+2} || \pi_{k+1}) \end{cases}$$

*so IQ-DP( $\alpha, \tau$ ) produces the same sequence of policies as a value-based version of Mirror Descent VI, MD-VI( $\alpha\tau, (1 - \alpha)\tau$ ) ([Vieillard et al., 2020a](#)).*

**Discussion.** The previous results shed a first light on the nature of the IQ method. Essentially, IQ-DP is a parametrization of a VI scheme regularized with both entropy and KL divergence, MD-VI( $\alpha\tau, (1 - \alpha)\tau$ ). This first highlights the role of the hyperparameters, as it shows the interaction between the two forms of regularization. The value of  $\alpha$  balances between those two: with  $\alpha = 0$ , IQ-DP reduces to a classic VI regularized with entropy; with  $\alpha = 1$  only the KL regularization will be taken into account. The value of  $\tau$  then controls

the amplitude of this regularization. In particular, in the limit  $\alpha = 0, \tau \rightarrow 0$ , we recover the standard VI algorithm. This result also justifies the soundness of IQ-DP. Indeed, this MD-VI scheme is known to converge to  $\pi_*^{(1-\alpha)\tau}$  the optimal policy of the regularized MDP (Vieillard et al., 2020a, Thm. 2) and this result readily applies to IQ<sup>3</sup>. Another consequence is that it links IQ to Advantage Learning (AL) (Bellemare et al., 2016). Indeed, AL is a limiting case of MD-VI when  $\alpha > 0$  and  $\tau \rightarrow 0$  (Vieillard et al., 2020b). Therefore, IQ also generalizes AL, and the  $\alpha$  parameter can be interpreted as the advantage coefficient. Finally, a key observation is that IQ performs KL regularization implicitly, the way it was introduced by Munchausen RL (Vieillard et al., 2020b), by augmenting the reward with the  $\alpha\tau \ln \pi_{k+1}$  term (Eq. (3)). This observation will have implications discussed next.

### 3.2 Error propagation result

Now, we are interested in understanding how errors introduced by the function approximation used propagate along iterations. At iteration  $k$  of IQ, denote  $\pi_{k+1}$  and  $V_k$  the target networks. In the approximate setting, we do not solve Eq. (3), but instead, we minimize  $\mathcal{L}(\theta, \phi)$  with stochastic gradient descent. This means that  $\pi_{k+2}$  and  $V_{k+1}$  are the result of this optimization, and thus the next target networks. The optimization process introduces errors, that come from many sources: partial optimization, function approximation (policy and value are approximated with neural networks), finite data, etc. We study the impact of these errors on the distance between the optimal  $Q$ -value of the MDP and the regularized  $Q$ -value of the current policy used by IQ,  $Q_{\pi_{k+1}}^{(1-\alpha)\tau}$ . We insist right away that  $Q_{\pi_{k+1}}^{(1-\alpha)\tau}$  is not the learned, implicit  $Q$ -value, but the actual  $Q$ -value of the policy computed by IQ in the regularized MDP. We have the following result concerning the error propagation.

**Theorem 2.** *Write  $\pi_{k+1}$  and  $V_k$  the  $k^{\text{th}}$  update of respectively the target policy and value networks. Consider the error at step  $k$ ,  $\epsilon_k \in \mathbb{R}^{S \times A}$ , as the difference between the ideal and the actual updates of IQ. Formally, we define the error as, for all  $k \geq 1$ ,*

$$\epsilon_k = \tau \ln \pi_{k+2} + V_{k+1} - (r + \alpha\tau \ln \pi_{k+1} + \gamma P V_k),$$

and the moving average of the errors as  $E_k = (1 - \alpha) \sum_{j=1}^k \alpha^{k-j} \epsilon_j$ . We have the following results for two different cases depending on the value of  $\alpha$ . Note that when  $\alpha < 1$ , we bound the distance to regularized optimal  $Q$ -value.

<sup>3</sup>Vieillard et al. (2020a) show this for  $Q$ -functions, but it can straightforwardly be extended to value functions.

1. *General case:  $0 < \alpha < 1$  and  $\tau > 0$ , entropy and KL regularization together:*

$$\|Q_*^{(1-\alpha)\tau} - Q_{\pi_k}^{(1-\alpha)\tau}\|_\infty \leq \frac{2}{(1-\gamma)^2} \left( (1-\gamma) \sum_{j=1}^k \gamma^{k-j} \|E_j\|_\infty \right) + o\left(\frac{1}{k}\right).$$

2. *Specific case  $\alpha = 1$ ,  $\tau > 0$ , use of KL regularization alone:*

$$\|Q_* - Q_{\pi_k}\|_\infty \leq \frac{2}{1-\gamma} \left\| \frac{1}{k} \sum_{j=1}^k \epsilon_j \right\|_\infty + O\left(\frac{1}{k}\right).$$

*Sketch of proof.* The full proof is provided in Appendix A.4. We build upon the connection we established between IQ-DP and a VI scheme regularized by both KL and entropy in Thm. 1. By injecting the proposed representation into the classic MD-VI scheme, we can build upon the analysis of Vieillard et al. (2020a, Thm. 1 and 2) to provide these results.  $\square$

**Impact of KL regularization.** The KL regularization term, and specifically in the MD-VI framework, is discussed extensively by Vieillard et al. (2020a), and we refer to them for in-depth analysis of the subject. We recall here the main interests of KL regularization, as illustrated by the bounds of Thm 2. In the second case, where it is the clearest (only KL is used), we observe a beneficial property of KL regularization: Averaging of errors. Indeed, in a classic non-regularized VI scheme (Scherrer et al., 2015), the error  $\|Q_* - Q_{\pi_\theta}\|$  would depend on a moving average of the norms of the errors  $(1-\gamma) \sum_{j=1}^k \gamma^{k-j} \|\epsilon_k\|_\infty$ , while with the KL it depends on the norm of the average of the errors  $(1/k) \|\sum_{j=1}^k \epsilon_k\|$ . In a simplified case where the errors would be i.i.d. and zero mean, this would allow convergence of approximate MD-VI, but not of approximate VI. In the case  $\alpha < 1$ , where we introduce entropy regularization, the impact is less obvious, but we still transform a sum of norm of errors into a sum of moving average of errors, which can help by reducing the underlying variance.

**Link to Munchausen RL.** As stated in the sketched proof, Thm. 2 is a consequence of (Vieillard et al., 2020a, Thm. 1 and 2). A crucial limitation of this work is that the analysis only applies when no errors are made in the greedy step. This is possible in a relatively simple setting, with tabular representation, or with a linear parametrization of the  $Q$ -function. However, in the general case with function approximation, exactly solving the optimization problem regularized by KL is not immediately possible: the solution of the greedy

step of MD-VI( $\alpha\tau, (1-\alpha)\tau$ ) is  $\pi_{k+2} \propto \exp(Q_{k+1}/\tau)\pi_k^\alpha$  (where  $Q_{k+1} = r + \gamma PV_k$ ), so computing it exactly would require remembering every  $\pi_j$  during the procedure, which is not feasible in practice. A workaround to this issue was introduced by Vieillard et al. (2020b) as Munchausen RL: the idea is to augment the reward by the log-policy, to implicitly define a KL regularization term, while reducing the greedy step to a softmax. As mentioned before, in small discrete action spaces, this allows to compute the greedy step exactly, but it is not the case in multidimensional or continuous action spaces, and thus Munchausen RL loses its interest in such domains. With IQ, we utilize the Munchausen idea to implicitly define the KL regularization; but with our parametrization, the exactness of the greedy step holds even for complex action spaces: recall that the parametrization defined in Eq. (1) enforces that the policy is a softmax of the (implicit)  $Q$ -value. Thus, IQ can be seen as an extension of Munchausen RL to multidimensional and continuous action spaces.

### 3.3 Link to the dueling architecture

Dueling Networks (DN) were introduced as a variation of the seminal Deep Q-Networks (DQN, Mnih et al. (2015)), and have been empirically proven to be efficient (for example by Hessel et al. (2018)). The idea is to represent the  $Q$ -value as the sum of a value and an advantage. In this setting, we work with a notion of advantage defined over  $Q$ -functions (as opposed to defining the advantage as a function of a policy). For any  $Q \in \mathbb{R}^{S \times A}$ , its advantage  $A_Q$  is defined as  $A_Q(s, a) = Q(s, a) - \max_{a' \in A} Q(s, a')$ . The advantage encodes a sub-optimality constraint: it has negative values and its maximum over actions is 0. Wang et al. (2016) propose to learn a  $Q$ -value by defining an advantage network  $F_\Theta$  and a value network  $V_\Phi$ , which in turn define a  $Q$ -value  $Q_{\Theta, \Phi}$  as

$$Q_{\Theta, \Phi}(s, a) = F_\Theta(s, a) - \max_{a' \in A} F_\Theta(s, a') + V_\Phi(s). \quad (4)$$

Subtracting the maximum over the actions ensures that  $F_\Theta$  indeed represents an advantage (in red). Note that dueling DQN was designed for discrete settings, where computing the maximum over actions is not an issue.

In IQ, we need a policy network that represents a distribution over the actions. There are several practical ways to represent the policy, discussed in Sec 4. For the sake of simplicity, let us for now assume that we are in a mono-dimensional discrete action space, and that we use a common scaled softmax representation. Specifically, our policy is represented by a neural network (eg. fully connected)  $F_\theta$ , that maps state-action pairs to logits  $F_\theta(s, a)$ . The policy is then defined as  $\pi_\theta(\cdot|s) = \text{softmax}(F_\theta(s, \cdot)/\tau)$ . Directly from the definition of the softmax, we observe that

$\tau \ln \pi_\theta(a|s) = F_\theta(s, a) - \tau \ln \sum_{a' \in A} \exp(F_\theta(s, a')/\tau)$ . The second term is a classic scaled logsumexp over the actions, a soft version of the maximum: when  $\tau \rightarrow 0$ , we have that  $\tau \ln \sum_{a'} \exp(F_\theta(s, a')/\tau) \rightarrow \max_{a'} F_\theta(s, a')$ . Within the IQ parametrization, we have

$$Q_{\theta, \phi}(s, a) = F_\theta(s, a) - \tau \ln \sum_{a' \in A} \exp \frac{F_\theta(s, a')}{\tau} + V_\phi(s),$$

where we highlighted in red the soft advantage component, analog to the advantage in Eq. (4), which makes a clear link between IQ and DN. In this case (scaled softmax representation), the IQ parametrization generalizes the dueling architecture, retrieved when  $\tau \rightarrow 0$  (and with an additional AL term whenever  $\alpha > 0$ , see Sec. 3). In practice, Wang et al. (2016) use a different parametrization of the advantage, replacing the maximum by a mean, defining  $Q_{\Theta, \Phi}(s, a) = A_\Theta(s, a) - |\mathcal{A}|^{-1} \sum_{a' \in A} A_\Theta(s, a') + V_\Phi(s)$ . We could use a similar trick and replace the logsumexp by a mean in our policy parametrization, but in our case this did not prove to be efficient in practice.

We showed how the log-policy represents a soft advantage. While this explicit its role in the learning procedure, it also raises questions about which representation would be the most suited for optimization.

## 4 PRACTICAL CONSIDERATIONS

We now describe key practical issues encountered when choosing a policy representation. The main one comes from the delegation of the representation power of the algorithm to the policy network. In a standard actor-critic algorithm – take SAC for example, where the policy is parametrized as a Gaussian distribution – the goal of the policy is mainly to track the maximizing action of the  $Q$ -value. Thus, estimation errors can cause the policy to choose sub-optimal actions, but the inductive bias caused by the Gaussian representation may not be a huge issue in practice, as long as the mean of the Gaussian policy is not too far from the maximizing action. In other words, the representation capacity of an algorithm such as SAC lies mainly in the representation capacity of its  $Q$ -network. In IQ, we have a parametrization of the policy that enforces it to be a softmax of an implicit  $Q$ -value. By doing this, we trade in estimation error – our greedy step is exact by construction – for representation power. Indeed, as our  $Q$ -value is not parametrized explicitly, but through the policy, the representation power of IQ is in its policy network, and a “simple” representation might not be enough anymore. For example, if we parameterized the policy as a Gaussian, this would amount to parametrize an advantage as a quadratic function of the action: this would drastically limit what the IQ policy could represent (in terms of soft-advantage).

**Multicategorical policies.** To address this issue, we turn to other, richer, distribution representations. In practice, we consider a multi-categorical discrete softmax distribution. Precisely, we are in the context of a multi-dimensional action space  $\mathcal{A}$  of dimension  $d$ , each dimension being a bounded interval. We discretize each dimension of the space uniformly in  $n$  values  $\delta_j$ , for  $0 \leq j \leq n - 1$ . It effectively defines a discrete action space  $\mathcal{A}' = \times_{j=1}^d \mathcal{A}_j$ , with  $\mathcal{A}_j = \{\delta_0, \dots, \delta_{n-1}\}$ . A multidimensional action is a vector  $a \in \mathcal{A}'$ , and we denote  $a^j$  the  $j^{\text{th}}$  component of the action  $a$ . Assuming independence between actions conditioned on states, a policy  $\pi_\theta$  can be factorized as the product of  $d$  marginal mono-dimensional policies  $\pi_\theta(a|s) = \prod_{j=1}^d \pi_\theta^j(a^j|s)$ . We represent each policy as the softmax of the output of a neural network  $F_\theta^j$ , and thus we get the full representation  $\pi_\theta(a|s) = \prod_{j=1}^d \text{softmax}(F_\theta^j(\cdot|s))(a^j)$ . The  $F_\theta^j$  functions can be represented as neural networks with a shared core, which only differ in the last layer. This type of multicategorical policy can represent any distribution (with  $n$  high enough) that does not encompass a dependency between the dimensions. The independence assumption is quite strong, and does not hold in general. From an advantage point of view, it assumes that the soft-advantage (*i.e.* the log-policy) can be linearly decomposed along the actions. While this limits the advantage representation, it is a much weaker constraint than parameterizing the advantage as a quadratic function of the action (which would be the case with a Gaussian policy). In practice, these types of policies have been experimented (Akkaya et al., 2019; Tang and Agrawal, 2020), and have proven to be efficient on continuous control tasks.

Even richer policy classes can be explored. To account for dependency between dimensions, one could envision auto-regressive multicategorical representations, used for example to parametrize a  $Q$ -value by Metz et al. (2017). Another approach is to use richer continuous distributions, such as normalizing flows (Rezende and Mohamed, 2015; Ward et al., 2019). In this work, we restrict ourselves to the multicategorical setting, which is sufficient to get satisfying results (Sec. 6), and we leave the other options for future work.

## 5 RELATED WORK

**Similar parametrizations.** Other algorithms use similar parametrization. First, Path Consistency Learning (PCL, (Nachum et al., 2017)) also parametrizes the  $Q$ -value as a sum of a log-policy and a value. Trust-PCL (Nachum et al., 2018), builds on PCL by adding a trust region constraint on the policy update, similar to our KL regularization term. A key difference with IQ is that (Trust-)PCL is a residual algorithm, while

IQ works around a fixed-point scheme. Shortly, Trust-PCL can be seen as a version of IQ without the target value network  $V_\phi$ . These entropy-regularized residual approaches are derived from the softmax temporal consistency principle, which allows to consider extensions to a specific form of multi-step learning (strongly relying on the residual aspect), but they also come with drawbacks, such as introducing a bias in the optimization when the environment is stochastic (Geist et al., 2017). Dai et al. (2018) proposed SBEED (Smooth Bellmann Error Embedding) to address this bias by replacing the residual minimization with a saddle-point problem. They provide an unbiased algorithm, but that ultimately requires to solve a min-max optimization, more complex than what we propose. Second, Quinoa (Degraeve et al., 2018) uses a similar loss to Trust-PCL and IQ (without reference to the former Trust-PCL), but does not propose any analysis, and is evaluated only on a few tasks. Third, Normalized Advantage Function (NAF, Gu et al. (2016)) is designed with similar principles. In NAF, a  $Q$ -value is parametrized as a value and an advantage, the latter being quadratic on the action. It matches the special case of IQ with a Gaussian policy, where we recover this quadratic parametrization. Finally, the action branching architecture (Tavakoli et al., 2018) was proposed to emulate Dueling Networks in the continuous case, by using a similar  $Q$ -value parametrization with a multicategorical policy representation. This architecture is recovered by IQ in the case  $\alpha = 0$  and  $\tau \rightarrow 0$ .

**Regularization.** Entropy and KL regularization are used by many other RL algorithms. Notably, from a dynamic programming perspective, IQ-DP(0,  $\tau$ ) performs the same update as SAC – an entropy regularized VI. This equivalence is however not true in the function approximation regime. Due to the empirical success of SAC and its link to IQ, it will be used as a baseline on continuous control tasks. Other algorithms also use KL regularization, *e.g.* Maximum a posteriori Policy Optimization (MPO, Abdolmaleki et al. (2018)). We refer to Vieillard et al. (2020a) for an exhaustive review of algorithms encompassed within the MD-VI framework.

## 6 EXPERIMENTS

**Environments and metrics.** We evaluate IQ first on the OpenAI Gym environment (Brockman et al., 2016). It consists of 5 locomotion tasks, with action spaces ranging from 3 (Hopper-v2) to 17 dimensions (Humanoid-v2). We use a rather long time horizon setting, evaluating our algorithm on 20M steps on each environment. We also provide result on the Adroit manipulation dataset (Rajeswaran et al., 2017), with a

similar setting of 20M environment steps. Adroit is a collection of 4 hand manipulation tasks, often used in an offline RL setting, but here we use it only as a direct RL benchmark. Out of these 4 tasks, we only consider 3 of them: we could not find any working algorithm (baseline or new) on the “relocate” task. To summarize the performance of an algorithm, we report the baseline-normalized score along iterations: It normalizes the score so that 0% corresponds to a random score, and 100% to a given baseline. It is defined for one task as  $\text{score} = \frac{\text{score}_{\text{algorithm}} - \text{score}_{\text{random}}}{\text{score}_{\text{baseline}} - \text{score}_{\text{random}}}$ , where the baseline is the best version of SAC on Mujoco and Adroit after 20M steps. We report aggregated results, showing the mean and median of these normalized scores along the tasks. Each score is reported as the average over 20 random seeds. For each experiment, the corresponding standard deviation is reported in Appx. B.3.

**IQ algorithms.** We implement IQ with the Acme (Hoffman et al., 2020) codebase. It defines two deep neural networks, a policy network  $\pi_\theta$  and a value network  $V_\phi$ . IQ interacts with the environment through  $\pi_\theta$ , and collect transitions that are stored in a FIFO replay buffer. At each interaction, IQ updates  $\theta$  and  $\phi$  by performing a step of stochastic gradient descent with Adam (Kingma and Ba, 2015) on  $\mathcal{L}_{\text{IQ}}$  (Eq. (2)). During each step, IQ updates a copy of the weights  $\theta$ ,  $\bar{\theta}$ , with a smooth update  $\bar{\theta} \leftarrow (1 - \lambda)\bar{\theta} + \lambda\theta$ , with  $\lambda \in (0, 1)$ . It tracks a similar copy  $\bar{\phi}$  of  $\phi$ . We keep almost all common hyperparameters (networks architecture,  $\lambda$ , etc.) the same as our main baseline, SAC. We only adjust the learning rate for two tasks, Humanoid and Walker, where we used a lower value: we found that IQ benefits from this, while for SAC we did not observe any improvement (we provide more details and complete results in Appx. B.3). Our value network has the same architecture as the SAC  $Q$ -networks except that the input size is only the state size (as it does not depend on the action). The policy network has the same architecture as the SAC policy network, and differs only by its output: IQ policy outputs a multicategorical policy (so  $n \cdot d$  values, where  $d$  is the dimensionality of the action space and  $n$  is the number of discrete action on each dimension), while SAC policy outputs 2  $d$ -dimensional vectors (mean and diagonal covariance matrix of a Gaussian). We use  $n = 11$  in our experiments. IQ relies on two other hyperparameters,  $\alpha$  and  $\tau$ . We selected a value of  $\tau$  per task suite:  $10^{-2}$  on Gym tasks and  $10^{-3}$  on Adroit. To make the distinction between the cases when  $\alpha = 0$  and  $\alpha > 0$ , we denote IQ( $\alpha > 0$ ) as M-IQ, for Munchausen-IQ, since it makes use of the Munchausen regularization term. For M-IQ, we found  $\alpha = 0.9$  to be the best performing value, which is consistent with the findings of Vieillard et al. (2020b). We report a empirical study of these

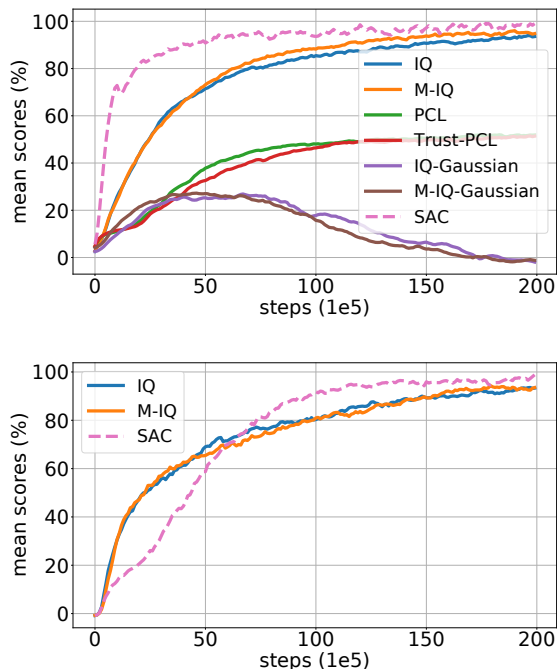


Figure 2: SAC-normalized mean scores on Gym (**top**) and Adroit (**bottom**).

parameters in the next part of this Section. Extended explanations are provided in Appx. B.2.

**Baselines.** On continuous control tasks, our main baseline is SAC, as it reaches state-of-the-art performance on Mujoco tasks. We compare to the version of SAC that uses an adaptive temperature for reference, but note that for IQ we keep a fixed temperature ( $\tau$ ) setting. To reach its best performance, SAC either uses a specific temperature value per task, or an adaptive scheme that controls the entropy of the policy. This method could be extended to multicategorical policies, but we leave this for future work, and for IQ we use the same value of  $\tau$  for all tasks of an environment. We use SAC with the default parameters from Haarnoja et al. (2018b) on Gym, and a specifically tuned version of SAC on Adroit. Remarkably, SAC and IQ work with similar hyperparameter ranges on both benchmarks. We only found that using a learning rate of  $3 \cdot 10^{-5}$  (instead of  $3 \cdot 10^{-4}$ ) gave better performance on Adroit. We also compare IQ to Trust-PCL. It is the closest algorithm to IQ, with a similar parametrization (SBEED also has this parametrization, but all the environments we consider here have deterministic dynamics, and thus there is no bias issue, see Sec. 5). To be fair, we compare to our version of Trust-PCL, which is essentially a residual version of IQ, where the target value network  $V_\phi$  is replaced by the online one. We use Trust-PCL with a fixed temperature, and we tuned this temperature to the environment. We found that Trust-PCL

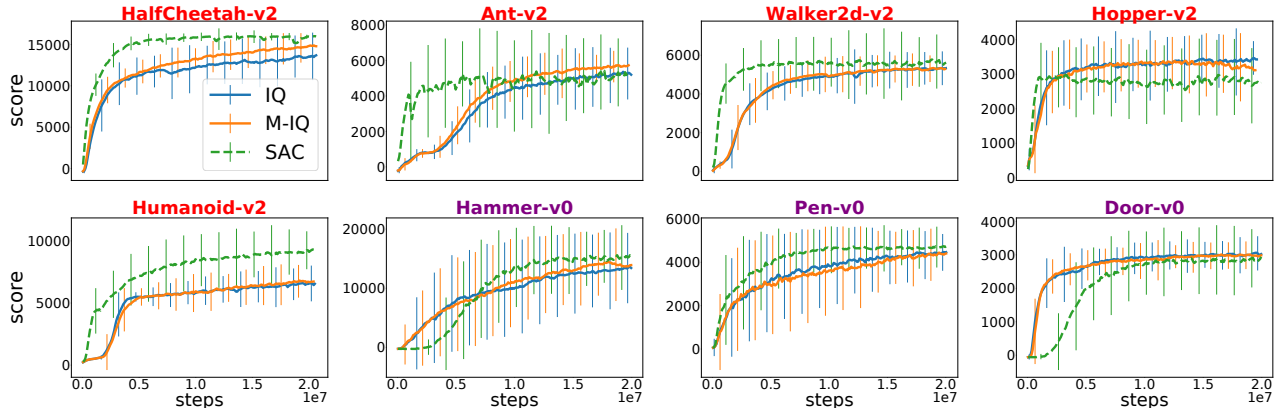


Figure 3: Scores on **Gym** and **Adroit**. Vertical bars denote  $\pm$  empirical standard deviation over 20 seeds.

reaches its best performance with significantly lower values of  $\tau$  compared to IQ. In the ablation (Fig. 2) we used  $\tau = 10^{-4}$  for PCL and Trust-PCL.

**Comparison to baselines.** We report aggregated results of IQ and M-IQ on Gym and Adroit in Fig. 2 (median scores can be found in Appx. B.3). IQ reaches competitive performance to SAC. It is less sample efficient on Gym (SAC reaches higher performance sooner), but faster on Adroit, and IQ reaches a close final performance on both environments. Detailed scores in Fig. 3 show how the performance varies across environments. Specifically, IQ outperforms SAC on 3 of the 8 considered environments: Ant, Hopper, and Door. On the 5 others, SAC performs better, but IQ and M-IQ still reach a reasonable performance. Moreover, on almost all environments, all mean performances are within the same confidence interval (Humanoid, and to a less extent HalfCheetah, being notable exceptions, in favor of SAC). The difference in performance between IQ and M-IQ represents the impact of the Munchausen term (*i.e.* KL regularization). It is never detrimental, and can even bring some improvement on Gym, but is clearly less useful than in discrete action. Indeed, M-DQN (Vieillard et al., 2020b) is designed with the same  $\alpha$  parameter but is tailored for discrete actions, and empirically clearly benefits from using  $\alpha > 0$ . We conjecture that this discrepancy comes from the inductive bias we introduce with the multicategorical policy, that could lessen the effect of the Munchausen trick.

**Influence of the hyperparameters.** (M)-IQ relies on three key hyperparameters: the temperature  $\tau$ , the Munchausen coefficient  $\alpha$ , and the number of bins  $n$ . The influence of  $\alpha$  can be understood as the difference between IQ and M-IQ, discussed above. For the two others parameters, we study their influence with two sweeps, for which we report learning curves in Appx B.3. IQ seems to be robust enough to the number of bins,

and  $n$  turns out to be rather easy to tune: values of  $n$  between 7 and 15 give similar results. On the other hand, we find that  $\tau$  needs to be selected carefully: while it helps learning, too high values of  $\tau$  can be detrimental to the performance, and it highlights that its optimal value might be dependant on the task. In the end,  $\tau$  has more influence on IQ than  $\alpha$  or  $n$ , and stands out to be the key element to tune within IQ; thus, adapting the automatic scheduling of the temperature from SAC looks like a promising research direction.

**Ablation study.** We perform an ablation on important components of IQ in Fig. 2. **(1)** We replace the target network by its online counterpart in Eq. (2), which gives us Trust-PCL (PCL is obtained by setting  $\alpha = 0$ ), a residual version of our method. IQ and M-IQ both outperform Trust-PCL and PCL on Mujoco. **(2)** We use a Gaussian parametrization of the policy instead of a multicategorical distribution. We observe on Fig. 2 that this causes the performance to drop drastically. This validates the considerations about the necessary complexity of the policy from Section 4.

## 7 CONCLUSION

We introduced IQ, a parametrization of a  $Q$ -value that mechanically preserves the softmax relation between a policy and an implicit  $Q$ -function. Building on this parametrization, we derived an off-policy algorithm, that learns a policy and a value by minimizing a single loss. We provided insightful analysis that justifies our algorithm and creates meaningful links with the literature (notably with the dueling networks architecture). Specifically, IQ performs entropy and implicit KL regularization on the policy. This kind of regularization was already used and analyzed in RL, but was limited by the difficulty of estimating the softmax of  $Q$ -function in continuous action settings. IQ ends this limitation by avoiding any approximation



in this softmax, effectively extending the analysis of this regularization. This parametrization comes at a cost: it shifts the representation capacity from the  $Q$ -network to the policy, which makes the use of Gaussian representation ineffective. We solved this issue by considering multicategorical policies, which allowed IQ to reach performance comparable to state-of-the-art methods on classic continuous control benchmarks. Yet, we envision that studying richer policy classes may result in even better performance. In the end, this work brings together theory and practice: IQ is a theory-consistent manner of implementing an algorithm based on regularized VI in continuous actions settings.

## References

- Abbas Abdolmaleki, Jost Tobias Springenberg, Yuval Tassa, Remi Munos, Nicolas Heess, and Martin Riedmiller. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations (ICLR)*, 2018. 2, 6
- Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, et al. Solving rubik’s cube with a robot hand. *arXiv preprint arXiv:1910.07113*, 2019. 6
- Mohammad G Azar, Mohammad Ghavamzadeh, Hilbert J Kappen, and Rémi Munos. Speedy Q-learning. In *Advances in Neural Information Processing System (NeurIPS)*, pages 2411–2419, 2011. 15
- Marc G Bellemare, Georg Ostrovski, Arthur Guez, Philip S Thomas, and Rémi Munos. Increasing the action gap: New operators for reinforcement learning. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2016. 4
- Dimitri P Bertsekas and John N Tsitsiklis. *Neuro dynamic programming*. Athena Scientific Belmont, MA, 1996. 1
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016. 2, 6
- Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. Sbed: Convergent reinforcement learning with nonlinear function approximation. In *International Conference on Machine Learning (ICML)*, pages 1125–1134. PMLR, 2018. 6
- Jonas Degraeve, Abbas Abdolmaleki, Jost Tobias Springenberg, Nicolas Heess, and Martin Riedmiller. Quinoa: a q-function you infer normalized over actions. *Deep RL Workshop at NeurIPS*, 2018. 6
- Thomas Degris, Martha White, and Richard S Sutton. Off-policy actor-critic. *International Conference on Machine Learning (ICML)*, 2012. 1
- Matthieu Geist, Bilal Piot, and Olivier Pietquin. Is the bellman residual a bad proxy? *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 6
- Matthieu Geist, Bruno Scherrer, and Olivier Pietquin. A Theory of Regularized Markov Decision Processes. In *International Conference on Machine Learning (ICML)*, 2019. 2, 3
- Shixiang Gu, Timothy Lillicrap, Ilya Sutskever, and Sergey Levine. Continuous deep q-learning with model-based acceleration. In *International Conference on Machine Learning (ICML)*, 2016. 6
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning (ICML)*, 2018a. 1, 2
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018b. 7
- Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2018. 5
- Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of convex analysis*. Springer Science & Business Media, 2004. 12
- Matt Hoffman, Bobak Shahriari, John Aslanides, Gabriel Barth-Maron, Feryal Behbahani, Tamara Norman, Abbas Abdolmaleki, Albin Cassirer, Fan Yang, Kate Baumli, Sarah Henderson, Alex Novikov, Sergio Gómez Colmenarejo, Serkan Cabi, Caglar Gulcehre, Tom Le Paine, Andrew Cowie, Ziyu Wang, Bilal Piot, and Nando de Freitas. Acme: A research framework for distributed reinforcement learning. *arXiv preprint arXiv:2006.00979*, 2020. URL <https://arxiv.org/abs/2006.00979>. 7, 15
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International Conference for Learning Representations (ICLR)*, 2015. 7
- Luke Metz, Julian Ibarz, Navdeep Jaitly, and James Davidson. Discrete sequential prediction of continuous actions for deep rl. *arXiv preprint arXiv:1705.05035*, 2017. 6

- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015. 5
- Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. Bridging the gap between value and policy based reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. 6
- Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. Trust-pcl: An off-policy trust region method for continuous control. *International Conference on Learning Representations (ICLR)*, 2018. 6
- Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017. 2, 6
- Daniilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning (ICML)*. PMLR, 2015. 6
- Bruno Scherrer, Mohammad Ghavamzadeh, Victor Gabillon, Boris Lesner, and Matthieu Geist. Approximate modified policy iteration and its application to the game of Tetris. *Journal of Machine Learning Research*, 16:1629–1676, 2015. 1, 4
- Yunhao Tang and Shipra Agrawal. Discretizing continuous action space for on-policy optimization. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 34, pages 5981–5988, 2020. 6
- Arash Tavakoli, Fabio Pardo, and Petar Kormushev. Action branching architectures for deep reinforcement learning. In *AAAI Conference on Artificial Intelligence (AAAI)*, volume 32, 2018. 6
- Nino Vieillard, Tadashi Kozuno, Bruno Scherrer, Olivier Pietquin, Rémi Munos, and Matthieu Geist. Leverage the average: an analysis of kl regularization in rl. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020a. 1, 3, 4, 6, 12, 14, 15
- Nino Vieillard, Olivier Pietquin, and Matthieu Geist. Munchausen reinforcement learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 2020b. 4, 5, 7, 8, 21
- Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International conference on machine learning (ICML)*, pages 1995–2003. PMLR, 2016. 2, 5
- Patrick Nadeem Ward, Ariella Smofsky, and Avishek Joey Bose. Improving exploration in soft-actor-critic with normalizing flows policies. *arXiv preprint arXiv:1906.02771*, 2019. 6
- Christopher JCH Watkins and Peter Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992. 3
- Brian D Ziebart. *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*. PhD thesis, University of Washington, 2010. 1

**Contents.** This Appendix is organized as follows

- Appendix A completes the theoretical claims of the paper
  - A.1 details the softmax consistency
  - A.2 gives background on KL and entropy regularization
  - A.3 proves Theorem 1 on equivalence between IQ-DP and regularized VI
  - A.4 proves Theorem 2 on error propagation in IQ-DP
  - A.5 details the relation between IQ and Munchausen DQN
- Appendix B provides details on experiments
  - B.1 gives information on our technical setup, regarding software and computing infrastructure
  - B.2 details the implementation of our method
  - B.3.1 gives a complete view of our experimental results
  - B.3.2 studies the influence of the hyperparameters

## A Analysis

This Appendix provides details and proofs on the IQ parameterization.

**Reminder on notations.** Throughout the Appendix, we use the following notations. Recall that we defined the action dot product as, for any  $u$  and  $v \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ ,

$$\langle u, v \rangle = \left( \sum_{a \in \mathcal{A}} u(s, a)v(s, a) \right)_s \in \mathbb{R}^{\mathcal{S}}.$$

We also slightly overwrite the  $+$  operator. Precisely, for any  $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ ,  $V \in \mathbb{R}^{\mathcal{S}}$ , we define  $Q + V \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  as

$$\forall (s, a) \in \mathcal{S} \times \mathcal{A}, (Q + V)(s, a) = Q(s, a) + V(s).$$

Write  $\mathbf{1} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  the constant function of value 1. For any  $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ , we define the softmax operator as

$$\text{softmax}(Q) = \frac{\exp(Q)}{\langle \mathbf{1}, \exp Q \rangle} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}},$$

where the fraction is overwritten as the addition operator, that is for any state-action pair  $(s, a)$ ,

$$\text{softmax}(Q)(a|s) = \frac{\exp Q(s, a)}{\sum_{a' \in \mathcal{A}} \exp Q(s, a')}.$$

### A.1 About the softmax consistency

First, we provide a detailed explanation of the consistency of the IQ parametrization. In Section 2, we claim that parametrizing a  $Q$ -value as  $Q = \tau \ln \pi + V$  enforces the relation  $\pi = \text{softmax}(Q/\tau)$ . This relation comes mechanically from the constraint that  $\pi$  is a distribution over actions. For the sake of precision, we provide a detailed proof of this claim as formalized in the following lemma.

**Lemma 1.** For any  $Q \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ ,  $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$ ,  $V \in \mathbb{R}^{\mathcal{S}}$ , we have

$$Q = \tau \ln \pi + V \Leftrightarrow \begin{cases} \pi = \text{softmax}\left(\frac{Q}{\tau}\right) \\ V = \tau \ln \langle \mathbf{1}, \exp \frac{Q}{\tau} \rangle \end{cases}. \quad (5)$$

*Proof.* Directly from the left hand side of Eq. (5), we have

$$\pi = \exp \frac{Q - V}{\tau}.$$

Since  $\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}$  ( $\pi$  is a distribution over the actions), we have

$$\begin{aligned} \langle \mathbf{1}, \pi \rangle = 1 &\Leftrightarrow \langle \mathbf{1}, \exp \frac{Q - V}{\tau} \rangle = 1 \\ &\Leftrightarrow \left( \exp \frac{-V}{\tau} \right) \langle \mathbf{1}, \exp \frac{Q}{\tau} \rangle = 1 \quad (V \text{ does not depend on the actions}) \\ &\Leftrightarrow V = \tau \ln \langle \mathbf{1}, \exp \frac{Q}{\tau} \rangle. \end{aligned}$$

And, for the policy, this gives

$$\pi = \exp \frac{Q - V}{\tau} = \exp \frac{Q - \tau \ln \langle \mathbf{1}, \exp \frac{Q}{\tau} \rangle}{\tau} = \frac{\exp(\frac{Q}{\tau})}{\langle \mathbf{1}, \exp \frac{Q}{\tau} \rangle} = \text{softmax} \frac{Q}{\tau}.$$

□

## A.2 Useful properties of KL-entropy-regularized optimization

The following proofs relies on some properties of the KL divergence and of the entropy. Consider the greedy step of MD-VI( $(1 - \alpha)\tau, \alpha\tau$ ), defined in Thm. 1

$$\pi_{k+2} = \operatorname{argmax}_{\pi \in \Delta_{\mathcal{A}}^{\mathcal{S}}} \langle \pi, r + \gamma PV_k \rangle + (1 - \alpha)\tau \mathcal{H}(\pi) - \alpha\tau \text{KL}(\pi || \pi_{k+1}). \quad (6)$$

Since the function  $\pi \rightarrow (1 - \alpha)\tau \mathcal{H}(\pi) - \alpha\tau \text{KL}(\pi || \pi_{k+1})$  is concave in  $\pi$ , this optimization problem can be tackled using properties of the Legendre-Fenchel transform (see for example [Hiriart-Urruty and Lemaréchal \(2004, Chap. E\)](#) for general definition and properties, and [Vieillard et al. \(2020a, Appx. A\)](#) for application to our setting). We quickly state two properties that are of interest for this work in the following Lemma.

**Lemma 2.** *Consider the optimization problem of Eq. (6). Write  $Q_{k+1} = r + \gamma PV_k$ , we have that*

$$\pi_{k+2} = \frac{\pi_{k+1}^{\alpha} \exp \frac{Q_{k+1}}{\alpha}}{\langle \mathbf{1}, \pi_{k+1}^{\alpha} \exp \frac{Q_{k+1}}{\alpha} \rangle}.$$

We also get a relation between the maximizer and the maximum

$$\langle \pi_{k+2}, r + \gamma PV_k \rangle + (1 - \alpha)\tau \mathcal{H}(\pi) - \alpha\tau \text{KL}(\pi || \pi_{k+1}) = \tau \ln \left\langle \pi_{k+1}^{\alpha}, \exp \frac{Q_{k+1}}{\alpha} \right\rangle.$$

*Proof.* See [Vieillard et al. \(2020a, Appx. A\)](#).

□

## A.3 Equivalence to MD-VI: proof of Theorem 1

We turn to the proof of Thm 1. This result formalizes an equivalence in the exact case between the IQ-DP scheme and a VI scheme regularized by entropy and KL divergence. Recall that we define the update of IQ-DP at step  $k$  as

$$\tau \ln \pi_{k+2} + V_{k+1} = r + \alpha\tau \ln \pi_{k+1} + \gamma PV_k \quad \text{IQ-DP}(\alpha, \tau). \quad (7)$$

Note that we are for now considering the scenario where this update is computed exactly. We will consider errors later, in Thm 2. Recall Thm. 1.

**Theorem 1.** *For any  $k \geq 1$ , let  $(\pi_{k+2}, V_{k+1})$  be the solution of IQ-DP at step  $k$ . We have that*

$$\begin{cases} \pi_{k+2} = \operatorname{argmax} \langle \pi, r + \gamma PV_k \rangle + (1 - \alpha)\tau \mathcal{H}(\pi) - \alpha\tau \text{KL}(\pi || \pi_{k+1}) \\ V_{k+1} = \langle \pi_{k+2}, r + \gamma PV_k \rangle + (1 - \alpha)\tau \mathcal{H}(\pi_{k+2}) - \alpha\tau \text{KL}(\pi_{k+2} || \pi_{k+1}) \end{cases}$$

so IQ-DP( $\alpha, \tau$ ) produces the same sequence of policies as a value-based version of Mirror Descent VI, MD-VI( $\alpha\tau, (1 - \alpha)\tau$ ) ([Vieillard et al., 2020a](#)).

*Proof.* Applying Lemma 1 to Eq. (7) gives

$$\begin{cases} \pi_{k+2} = \text{softmax} \frac{r + \alpha\tau \ln \pi_{k+1} + \gamma PV_k}{\tau} \\ V_{k+1} = \tau \ln \left\langle \mathbf{1}, \exp \frac{r + \alpha\tau \ln \pi_{k+1} + \gamma PV_k}{\tau} \right\rangle. \end{cases}$$

For the policy, we have

$$\pi_{k+2} = \frac{\exp(\alpha \ln \pi_{k+1}) \exp \frac{r + \gamma PV_k}{\alpha}}{\left\langle \mathbf{1}, \exp(\alpha \ln \pi_{k+1}) \exp \frac{r + \gamma PV_k}{\alpha} \right\rangle} = \frac{\pi_{k+1}^\alpha \exp \frac{r + \gamma PV_k}{\alpha}}{\left\langle \mathbf{1}, \pi_{k+1}^\alpha \exp \frac{r + \gamma PV_k}{\alpha} \right\rangle},$$

and as direct consequence of Lemma 2

$$\pi_{k+2} = \text{argmax} \langle \pi, r + \gamma PV_k \rangle + (1 - \alpha)\tau \mathcal{H}(\pi) - \alpha\tau \text{KL}(\pi || \pi_{k+1}).$$

For the value, we have:

$$V_{k+1} = \tau \ln \left\langle \mathbf{1}, \exp(\alpha \ln \pi_{k+1}) \exp \frac{r + \gamma PV_k}{\tau} \right\rangle = \tau \ln \left\langle \pi_{k+1}^\alpha, \exp \frac{r + \gamma PV_k}{\tau} \right\rangle,$$

and again applying Lemma 2 gives

$$V_{k+1} = \langle \pi_{k+2}, r + \gamma V_k \rangle + (1 - \alpha)\tau \mathcal{H}(\pi_{k+2}) - \alpha\tau \text{KL}(\pi_{k+2} || \pi_{k+1}).$$

□

#### A.4 Error propagation: proof of Theorem 2

Now we turn to the proof of Thm 2. This theorem handles the IQ-DP scheme in the approximate case, when errors are made during the iterations. The considered scheme is

$$\tau \ln \pi_{k+2} + V_{k+1} = r + \alpha\tau \ln \pi_{k+1} + \gamma PV_k + \epsilon_{k+1}. \quad (8)$$

Recall Thm. 2.

**Theorem 2.** Write  $\pi_{k+1}$  and  $V_k$  the  $k^{\text{th}}$  update of respectively the target policy and value networks. Consider the error at step  $k$ ,  $\epsilon_k \in \mathbb{R}^{S \times A}$ , as the difference between the ideal and the actual updates of IQ. Formally, we define the error as, for all  $k \geq 1$ ,

$$\epsilon_{k+1} = \tau \ln \pi_{k+2} + V_{k+1} - (r + \alpha\tau \ln \pi_{k+1} + \gamma PV_k),$$

and the moving average of the errors as

$$E_k = (1 - \alpha) \sum_{j=1}^k \alpha^{k-j} \epsilon_j.$$

We have the following results for two different cases depending on the value of  $\alpha$ . Note that when  $\alpha < 1$ , we bound the distance to regularized optimal  $Q$ -value.

1. General case:  $0 < \alpha < 1$  and  $\tau > 0$ , entropy and KL regularization together:

$$\|Q_*^{(1-\alpha)\tau} - Q_{\pi_k}^{(1-\alpha)\tau}\|_\infty \leq \frac{2}{(1-\gamma)^2} \left( (1-\gamma) \sum_{j=1}^k \gamma^{k-j} \|E_j\|_\infty \right) + o\left(\frac{1}{k}\right).$$

2. Specific case  $\alpha = 1$ ,  $\tau > 0$ , use of KL regularization alone:

$$\|Q_* - Q_{\pi_k}\|_\infty \leq \frac{2}{1-\gamma} \left\| \frac{1}{k} \sum_{j=1}^k \epsilon_j \right\|_\infty + O\left(\frac{1}{k}\right).$$

*Proof.* To prove this error propagation result, we first show an extension of Thm. 1, that links Approximate IQ-DP with a  $Q$ -value based version of MD-VI. This new equivalence makes IQ-DP corresponds exactly to a scheme that is extensively analyzed by [Vieillard et al. \(2020a\)](#). Then our result can be derived as a consequence of ([Vieillard et al., 2020a](#), Thm 1) and ([Vieillard et al., 2020a](#), Thm 2).

Define a (KL-regularized) implicit  $Q$ -value as

$$Q_k = \tau \ln \pi_{k+1} - \alpha \tau \ln \pi_k + V_k,$$

so that now, the IQ-DP update (Eq. (8)) can be written

$$Q_{k+1} = r + \gamma P V_k + \epsilon_{k+1}. \quad (9)$$

We then use same method that for the proof of Thm. 1. Specifically, applying Lemma 1 to the definition of  $Q_k$  gives for the policy

$$\begin{aligned} \pi_{k+1} &= \text{softmax} \left( \frac{Q_k + \alpha \tau \ln \pi_k}{\tau} \right) \quad (\text{Lemma 1}) \\ &= \frac{\pi_k^\alpha \exp \frac{Q_k}{\alpha}}{\left\langle \mathbf{1}, \pi_k^\alpha \exp \frac{Q_k}{\alpha} \right\rangle} \\ \Leftrightarrow \pi_{k+1} &= \text{argmax} \langle \pi, Q_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi) - \alpha \tau \text{KL}(\pi || \pi_k). \quad (\text{Lemma 2}) \end{aligned}$$

For the value, we have from Lemma 1 on  $Q_k$

$$V_k = \tau \ln \left\langle \mathbf{1}, \exp \frac{Q_k + \alpha \tau \ln \pi_k}{\tau} \right\rangle = \tau \ln \left\langle \pi_k^\alpha, \exp \frac{Q_k}{\tau} \right\rangle,$$

then, using Lemma 2, and the fact that  $\pi_{k+1} = \text{softmax} \frac{Q_k}{\tau}$ , we have

$$V_k = \langle \pi_{k+1}, Q_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi_{k+1}) - \alpha \tau \text{KL}(\pi_{k+1} || \pi_k).$$

Injecting this in Eq. (9) gives

$$Q_{k+1} = r + \gamma P (\langle \pi_{k+1}, Q_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi_{k+1}) - \alpha \tau \text{KL}(\pi_{k+1} || \pi_k)).$$

Thus, we have proved the following equivalence between DP schemes

$$\begin{aligned} \tau \ln \pi_{k+2} + V_{k+1} &= r + \alpha \tau \ln \pi_{k+1} + \gamma P V_k + \epsilon_{k+1} \\ &\Updownarrow \\ \begin{cases} \pi_{k+1} = \text{argmax} \langle \pi, Q_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi) - \alpha \tau \text{KL}(\pi || \pi_k) \\ Q_{k+1} = r + \gamma P (\langle \pi_{k+1}, Q_k \rangle + (1 - \alpha) \tau \mathcal{H}(\pi_{k+1}) - \alpha \tau \text{KL}(\pi_{k+1} || \pi_k)) + \epsilon_{k+1}, \end{cases} \end{aligned} \quad (10)$$

with

$$Q_k = \tau \ln \pi_{k+1} - \alpha \tau \ln \pi_k + V_k.$$

The above scheme in Eq. (10) is exactly the MD-VI scheme studied by [Vieillard et al. \(2020a\)](#), where they define  $\beta = \alpha$  and  $\lambda = \alpha \tau$ . We now use their analysis of MD-VI to apply their result to IQ-DP, building on the equivalence between the schemes. Note that transferring this type of analysis between equivalent formulations of DP schemes is justified because the equivalences exist in terms of *policies*. Indeed, IQ-DP and MD-VI compute different ( $Q$ -)values, but produce *identical series of policies*. Since ([Vieillard et al., 2020a](#), Thm 1) and ([Vieillard et al., 2020a](#), Thm. 2) bound the distance between the optimal (regularized)  $Q$ -value and the actual (regularized)  $Q$ -values of the computed policy, the equivalence in terms of policies is sufficient to apply these theorems to IQ-DP. Specifically, ([Vieillard et al., 2020a](#), Thm 1) applied to the formulation of IQ in Eq. (10) proves point 1 of Thm. 2, that is the case where  $\alpha = 0$ . The second part is proven by applying ([Vieillard et al., 2020a](#), Thm 2) to this same formulation.  $\square$

**Remark on sample complexity.** In the case of a tabular representation (no function approximation needed), and with access to a generative model, we could derive a sample complexity bound from Theorem 2. Indeed, we can follow a similar analysis to the one of Speedy Q-learning by Azar et al. (2011) by combining our results with Azuma-Hoeffding concentration inequalities; it would show that IQ benefits from the same sample complexity as tabular MD-VI (Vieillard et al., 2020a), that is  $\mathcal{O}\left(\frac{|S||A|}{\epsilon^2(1-\gamma)^4}\right)$  samples to get an  $\epsilon$ -optimal policy.

### A.5 IQ and Munchausen DQN

We claim in Section 3 that IQ is a form of Munchausen algorithm, specifically Munchausen-DQN (M-DQN). Here, we clarify this link. Note that all of the information below is contained in Appx. A.3 and Appx. A.4. The point of this section is to re-write it using notations used to defined IQ as a deep RL agent, notations consistent with how M-DQN is defined.

IQ optimizes a policy  $\pi_\theta$  and a value  $V_\phi$  by minimizing a loss  $\mathcal{L}_{\text{IQ}}$  (Eq. (2)). Recall that IQ implicitly defines a Q-function as  $Q_{\theta,\phi} = \tau \ln \pi_\theta + V_\phi$ . Identifying this in  $\mathcal{L}_{\text{IQ}}$  makes the connection between Munchausen RL and IQ completely clear. Indeed, the loss can be written as

$$\hat{\mathbb{E}} \left[ \left( r_t + \alpha \tau \ln \pi_{\bar{\theta}}(a_t|s_t) + \underbrace{\frac{\gamma V_{\bar{\phi}}(s_{t+1})}{\tau \ln \sum_a \exp \frac{Q_{\bar{\theta},\bar{\phi}}(s_{t+1},a)}}}_{Q_{\bar{\theta},\bar{\phi}}} - \underbrace{\tau \ln \pi_{\bar{\theta}}(a_t|s_t) - V_{\bar{\phi}}(s_t)}_{Q_{\theta,\phi}} \right)^2 \right],$$

and since we have (Lemma 2, and using the fact that  $\pi_{\bar{\theta}} = \text{softmax}(Q_{\bar{\theta},\bar{\phi}}/\tau)$ )

$$\tau \ln \sum_a \exp \frac{Q_{\bar{\theta},\bar{\phi}}(s,a)}{\tau} = \sum_a \pi_{\bar{\theta}}(a|s) (Q_{\bar{\theta},\bar{\phi}}(s,a) - \tau \ln \pi_{\bar{\theta},\bar{\phi}}(a|s)),$$

we get that the loss is

$$\hat{\mathbb{E}} \left[ \left( r_t + \alpha \tau \ln \pi_{\bar{\theta}}(a_t|s_t) + \sum_a \pi_{\bar{\theta}}(a|s_{t+1}) (Q_{\bar{\theta},\bar{\phi}}(s_{t+1},a) - \tau \ln \pi_{\bar{\theta}}(a|s_{t+1})) - Q_{\theta,\phi}(s_t, a_t) \right)^2 \right],$$

which is exactly the Munchausen-DQN loss on  $Q_{\theta,\phi}$ . Thus, in a mono-dimensional action setting (classic discrete control problems for example), IQ can really be seen as a re-parameterized version of M-DQN.

## B Additional material on experiments

This Appendix provides additional detail on experiments, along with complete empirical results.

### B.1 General information on experiments

**Used assets.** IQ is implemented on the Acme library (Hoffman et al., 2020), distributed as open-source code under the Apache License (2.0).

**Compute resources.** Experiments were run on TPUv2. One TPU is used for a single run, with one random seed. To produce the main results (without the sweeps over parameters), we computed 780 single runs. One of this run on a TPUv2 takes from 3 to 10 hours depending on the environment (the larger the action space, the longer the run).

### B.2 Details on algorithms

**On the relation between  $\alpha$  and  $\tau$ .** The equivalence result of Theorem 1 explains the role and the relation between  $\tau$  and  $\alpha$ . In particular, it shows that IQ-DP( $\alpha, \tau$ ) performs a VI scheme in an entropy-regularized MDP (or in a max-entropy setting) where the temperature is not  $\tau$ , but  $(1 - \alpha)\tau$ . Indeed, in this framework, the  $\alpha$  parameter balances between two forms of regularization: with  $\alpha = 0$ , IQ-DP is only regularized with entropy,

but with  $\alpha > 0$ , IQ-DP is regularized with both entropy and KL. Thus, IQ-DP modifies implicitly the intrinsic temperature of the MDP it is optimizing for. To account for this discrepancy, every time we evaluate IQ with  $\alpha > 0$  (that is, M-IQ), we report scores using  $\tau/(1 - \alpha)$ , and not  $\tau$ . For example, on Gym, we used a temperature of 0.01 for IQ, and thus 0.1 for M-IQ (since, in our experiments, we took  $\alpha = 0.9$ ).

**Discretization.** We used IQ with policies that discretize the action space evenly. Here, we provide a precise definition for our discretization method. Consider a multi-dimensional action space  $\mathcal{A}$  of dimension  $d$ , each dimension being a bounded interval  $[a_{\min}, a_{\max}]$ , such that  $\mathcal{A} = [a_{\min}, a_{\max}]^d$ . We discretize each dimension of the space uniformly in  $n$  values  $\delta_j$ , for  $0 \leq j \leq n - 1$ . The bins values are defined as

$$\delta_0 = a_{\min} + \frac{a_{\max} - a_{\min}}{2n},$$

and, for each  $j \in \{1, \dots, n - 1\}$ ,

$$\delta_j = \delta_0 + j \frac{a_{\max} - a_{\min}}{n}.$$

It effectively defines a discrete action space

$$\mathcal{A}' = \bigtimes_{j=1}^d \mathcal{A}_j, \quad \text{with } \mathcal{A}_j = \{\delta_0, \dots, \delta_{n-1}\}.$$

We use  $n = 11$  in all of our experiments. The values of  $d$ ,  $a_{\min}$  and  $a_{\max}$  depend on the environments specifications.

**Evaluation setting.** We evaluate our algorithms on Mujoco environments from OpenAI Gym and from the Adroit manipulation tasks. On each environment, we track performance for 20M environment steps. Every 10k environment steps, we stop learning, and we evaluate our algorithm by reporting the average undiscounted return over 10 episodes. We use deterministic evaluation, meaning that, at evaluation time, the algorithms interact by choosing the expected value of the policy in one state, not by sampling from this policy (sampling is used during training).

**Pseudocode.** We provide a pseudocode of IQ in Algorithm 1. This pseudocode describes a general learning procedure that is followed by all agents. Replacing the IQ loss in Algorithm 1 by its residual version will give the pseudocode for PCL, and replacing it by the actor and critic losses of SAC will give the pseudocode for this method.

---

**Algorithm 1** Implicit Q-values

---

**Require:**  $T \in \mathbb{N}^*$  the number of environment steps,  $\lambda \in (0, 1)$  the update coefficient,  $\gamma \in [0, 1)$  the discount factor,  $\tau \in (0, 1)$  the entropy temperature,  $\alpha \in [0, 1)$  the implicit KL term, and hyperparameters detailed in Table 1.

Initialize  $\theta, \phi$  at random

$\mathcal{B} = \{\}$

$\bar{\theta} = \theta$

$\bar{\phi} = \phi$

**for**  $t = 1$  **to**  $T$  **do**

Collect a transition  $b = (s_t, a_t, r_t, s_{t+1})$  from  $\pi_\theta$

$\mathcal{B} \leftarrow \mathcal{B} \cup \{b\}$

On a random batch of transitions  $B_t \subset \mathcal{B}$ , update  $(\theta, \phi)$  with one step of SGD on

$$\hat{\mathbb{E}}_{(s_t, a_t, r_t, s_{t+1}) \sim B_t} \left[ (r_t + \alpha \tau \ln \pi_{\bar{\theta}}(a_t | s_t) + \gamma V_{\bar{\phi}}(s_{t+1}) - \tau \ln \pi_\theta(a_t | s_t) - V_\phi(s_t))^2 \right],$$

$\bar{\theta} \leftarrow \lambda \bar{\theta} + (1 - \lambda) \theta$

$\bar{\phi} \leftarrow \lambda \bar{\phi} + (1 - \lambda) \phi$

**end for**

**return**  $\pi_\theta$

---



**HyperParameters.** We Provide the hyperparameters used for our experiments in Tab. 1. If a parameter is under “common parameters”, then it was used for all algorithms. We denote  $FC_n$  a fully connected layer with an output of  $n$  neurons. Recall that  $d$  is the dimension of the action space, and  $n$  is the number of bins we discretize each dimension into.

Table 1: Parameters used for algorithms and ablations.

Parameter	Value
Common parameters	
$\lambda$ (update coefficient)	0.05
$\gamma$ (discount)	0.99
$ \mathcal{B} $ (replay buffer size)	$10^6$
$ B_t $ (batch size)	256
activations	Relu
optimizer	Adam
learning rate	$3 \cdot 10^{-4}$
IQ specific parameters	
$\tau$ (entropy temperature)	0.01 on Gym, 0.001 on Adroit
$\alpha$ (implicit KL term)	0.9
$n$ (number of bins for the discretization)	11
$\pi$ -network	(input: state) FC 512 – FC 512 – FC $nd$
$V$ -network structure	(input: state) FC 512 – FC 512 – FC 1
(Trust)-PCL specific parameters	
$\tau$ (entropy temperature)	$1 \cdot 10^{-4}$ on Gym
$\pi$ -network and $V$ -network structures	idem as IQ
SAC specific parameters	
$\pi$ -network structure	(input: state) FC 512 – FC 512 – FC $2d$
$Q$ -network structure	(input: state and action) FC 512 – FC 512 – FC 1
learning rate	$3 \cdot 10^{-4}$ on Gym, $3 \cdot 10^{-5}$ on Adroit

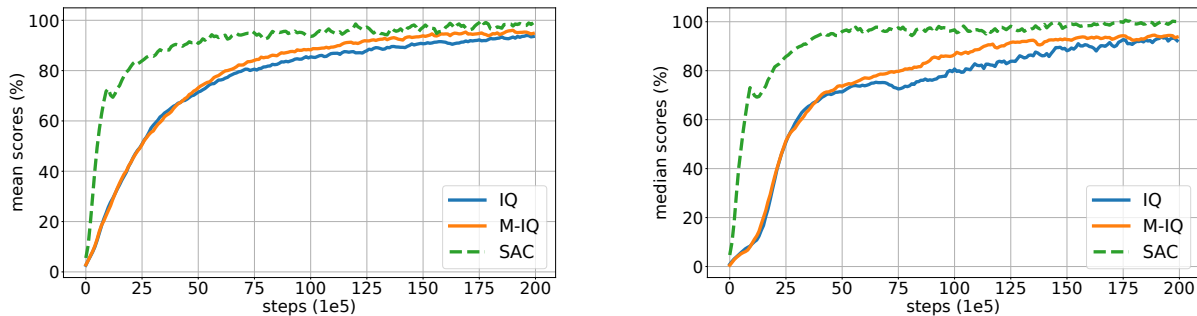


Figure 4: SAC-normalized aggregated scores on Gym environments. *Best parameters*: IQ uses a different learning rate for Humanoid-v2 and Walker2d-v2. **Left**: Mean scores. **Right**: Median scores.

### B.3 Additional results

This Appendix section provides complete description of the scores and ablations in B.3.1, along with a study of some hyperparameters in B.3.2.

#### B.3.1 Complete scores and ablations

Here, we provide detailed results in a more readable fashion in Fig. 4 and 5. We also give more details on the aggregated scores (median scores, explanations on the hyperparameter selection), and provide a more extensive analysis of our experiments.

**Aggregated scores.** Aggregated scores on Gym and Adroit environments are reported in Figures 9 and 5. We see from these results that (M)-IQ reaches performance close to SAC, but is still slightly below. On Gym, IQ is slower (as in less sample efficient) than SAC, but is faster on Adroit. One important thing to notice is that we compare to the best version of SAC we found, that uses an adaptive temperature, while we use the same fixed temperature for all of the environments. Experimentally, it appears that IQ is more sensitive to hyperparameters than SAC<sup>4</sup>. Notably, on Humanoid-v2 and Walker2d-v2, we show in Figure 6 that IQ can reach a much more competitive score by using a specific hyperparameter per environment. To reflect this, we use a learning rate of  $3 \cdot 10^{-5}$  instead of  $3 \cdot 10^{-4}$  for Humanoid and Walker, while still comparing to the best version of SAC we found (we also tested SAC with those parameters variations, but it did not improve the performance). This is reported as *best parameters* results in Fig. 4; and for completeness, we report score using *unique hyperparameters* in Fig. 9. On Adroit, we did not observe such variations, and thus we report scores with a single set of hyperparameters.

**Detailed scores.** We report detailed scores on Gym environments in Figure 7 and on Adroit in Figure 8. On this figures, the thick line corresponds to the average over 20 seeds, and the error bars represents  $+/-$  the empirical standard deviation over those seeds. It shows how the performance varies across environments. Specifically, IQ outperforms SAC on 3 of the 8 considered environments: Ant, Hopper, and Door. On the 5 others, SAC performs better, but IQ and M-IQ still reach a reasonable performance. Moreover, on almost all environments, all mean performances are within the same confidence interval (Humanoid, and to a less extent HalfCheetah, being notable exceptions, in favor of SAC).

**Ablations.** We report ablation scores, averaged over 20 seeds, comparing IQ, M-IQ, PCL, Trust-PCL and (M)-IQ-Gaussian in Fig. 9. PCL (and Trust-PCL) are obtained by replacing  $V_{\bar{\phi}}$  by  $V_{\phi}$  in  $\mathcal{L}_{IQ}$ . (M)-IQ Gaussian is obtained by parametrizing the policy as a diagonal normal distribution over the action. Results confirm that the fixed-point approach is beneficial wrt the residual one in this case, since IQ outperforms PCL and trust-PCL by a

<sup>4</sup>However, we note that SAC, when used with a fixed temperature, is sensitive to the choice of this parameter. It appears that IQ is less sensitive to this (we use the same temperature for all tasks), but a bit more to the learning rate. We think this may be alleviated by adopting other policy representations, such as normalizing flows, or by designing adaptive schemes. We left this for future works.

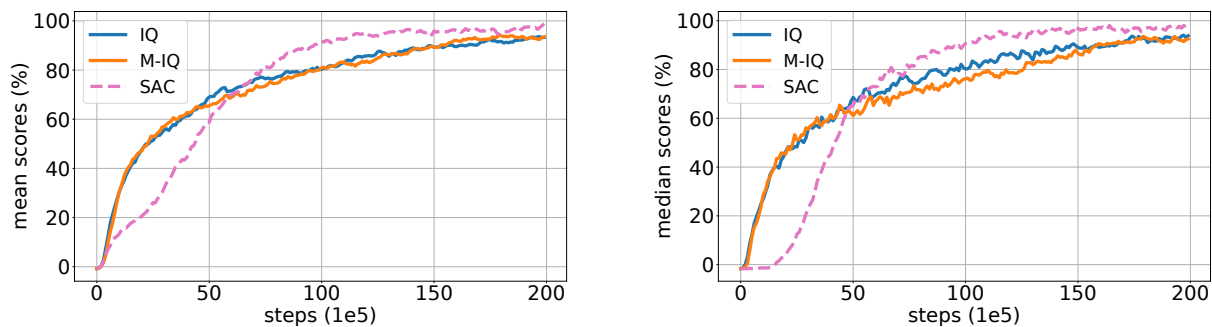


Figure 5: SAC-Normalized aggregated scores on Adroit. **Left:** Mean scores. **Right:** Median scores.

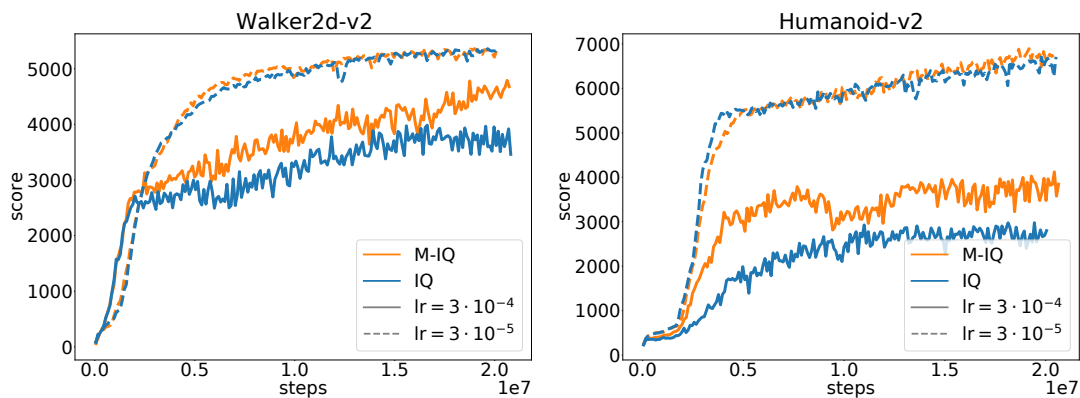


Figure 6: Comparison of two values of the learning rate ( $lr$ ) on IQ and M-IQ, on two environments. Each line corresponds to the average score over 20 seeds.

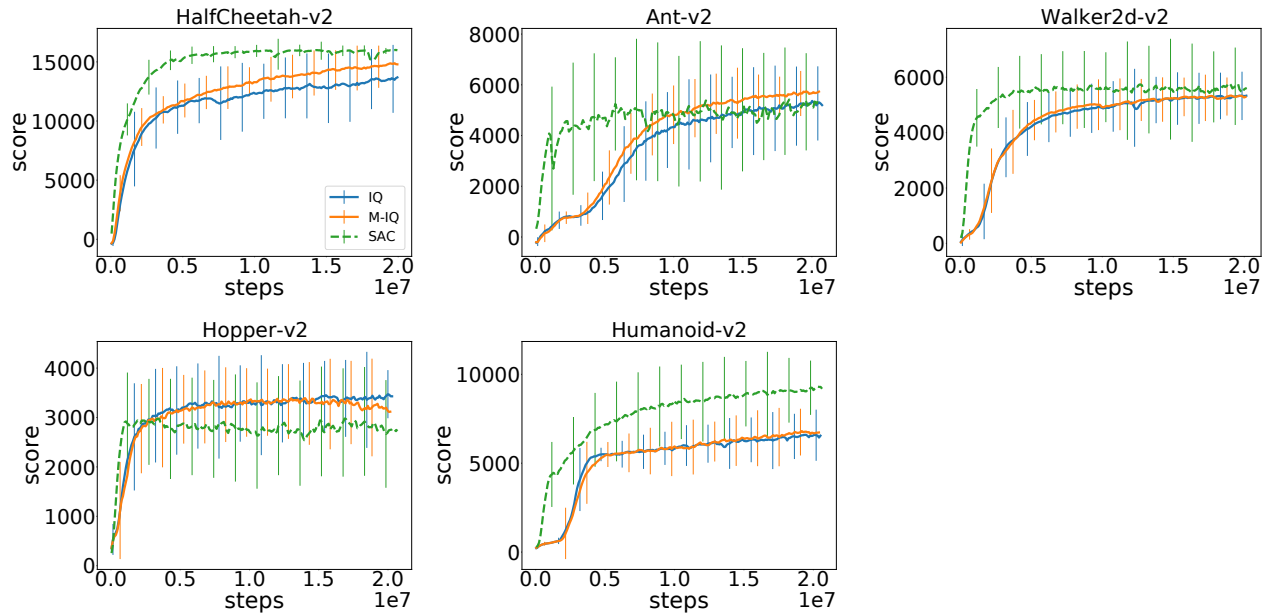


Figure 7: All individual scores on Gym. The vertical bars denote the empirical standard deviation over 20 seeds. *Best parameters:* IQ uses a different learning rate for Humanoid-v2 and Walker2d-v2.

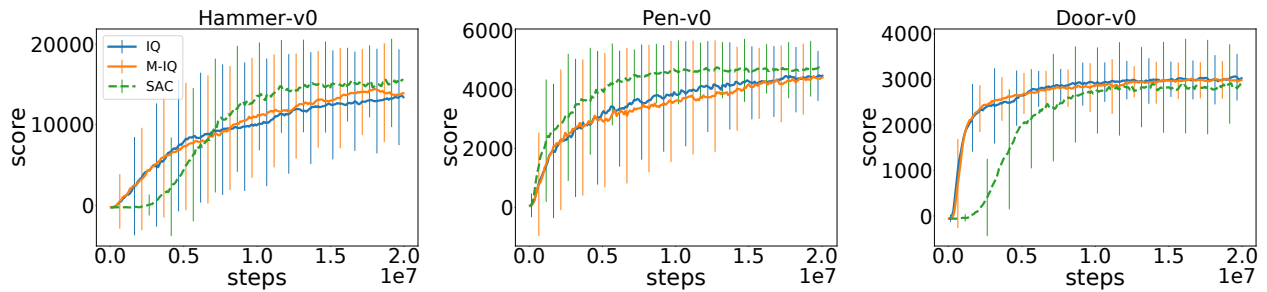


Figure 8: All scores on Adroit. The vertical bars denote the empirical standard deviation over 20 seeds.

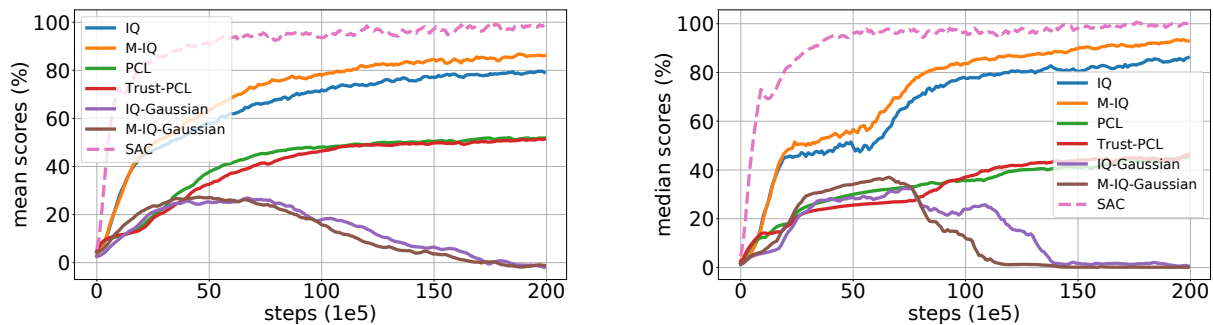


Figure 9: SAC-normalized ablation scores on Gym. *Unique parameters*: all algorithms use a single set of hyperparameters for the 5 environments.

clear margin in median and mean<sup>5</sup>. They also validate that the Gaussian parametrization is too limiting in IQ, since IQ-Gaussian and M-IQ-Gaussian both totally collapse after 6 to 7 millions environment steps. Additionally, this also highlights that the time frame usually considered on some of this environments (3M, or even 1M steps) may be too short to actually observe the asymptotic behavior of many agents. Here, this time frame would have prevented us to observe this collapsing phenomenon on the Gaussian IQ.

### B.3.2 Influence of hyperparameters

In addition to the "deep" hyperparameters that control the function approximation (layer sizes, activations, learning rate, ..), and that we take from our baselines, (M)-IQ relies on three key hyperparameters. Those are the Munchausen coefficient  $\alpha$ , the temperature  $\tau$ , and the number of bins  $n$ . The influence of  $\alpha$  is represented by the difference between M-IQ and IQ. Here, we study the influence of the two others hyperparameters,  $\tau$  and  $n$ .

**Influence of the temperature ( $\tau$ ).** We provide results for 3 values of  $\tau$  and 2 values of  $\alpha$ . We report scores for IQ with  $\alpha = 0$  with different values of  $\tau$  in Fig. 10 and the same experiments on M-IQ (with  $\alpha = 0.9$ ) in Fig. 11. We can draw two observations from these results. **(1)** They show that  $\tau$  needs to be selected carefully: it helps learning, but too high values can be highly detrimental to the performance. The temperature also depends on the task. While we could find a value that give reasonable results on all the considered environment, this means that adapting the automatic scheduling of the temperature used in SAC could be a promising direction for IQ. **(2)** The other observation is that clearly,  $\tau$  has a much stronger influence on IQ than  $\alpha$ . This is a key empirical difference regarding the performance of M-DQN (Vieillard et al., 2020b), that has the same parameters, but is evaluated on discrete actions settings. In these settings, the  $\alpha$  parameters is shown to have a crucial importance in terms of empirical results: M-DQN with  $\alpha = 0.9$  largely outperforms M-DQN with  $\alpha = 0$  on the Atari benchmark. While this term still has effect in IQ on some tasks, it is empirically less useful, even though it is never detrimental. We conjecture that this discrepancy comes from the inductive bias we introduce with the multicategorical policy, that could lessen the effect of the Munchausen trick.

**Influence of the number of bins.** In Fig. 12, we show how the number of bins impact the performance of IQ on the Gym environments. The results indicate that this parameter is fairly easy to tune, as IQ is not too sensitive to its variations. A remarkable observation is that with only 3 bins, IQ can achieve reasonable performance on complex environments (see the results on Humanoid for example). This turns out to be a beneficial feature of the multicategorical policies. Indeed, the limiting factor of this parametrization, is that we need a potentially large number of bins in total, if the dimensions of the environment grows: the output size of the policy network is  $d * n$ , with  $d$  the dimension of the action space and  $n$  the number of bins. If we can keep the number of bins close to a minimal value (as it seems to be the case here), this is promising for potential use of IQ on larger action spaces.

<sup>5</sup>We note that PCL was introduced only for discrete actions, and that Trust-PCL was proposed with a Gaussian policy, on a much shorter time frame for training, and with many additional tricks.

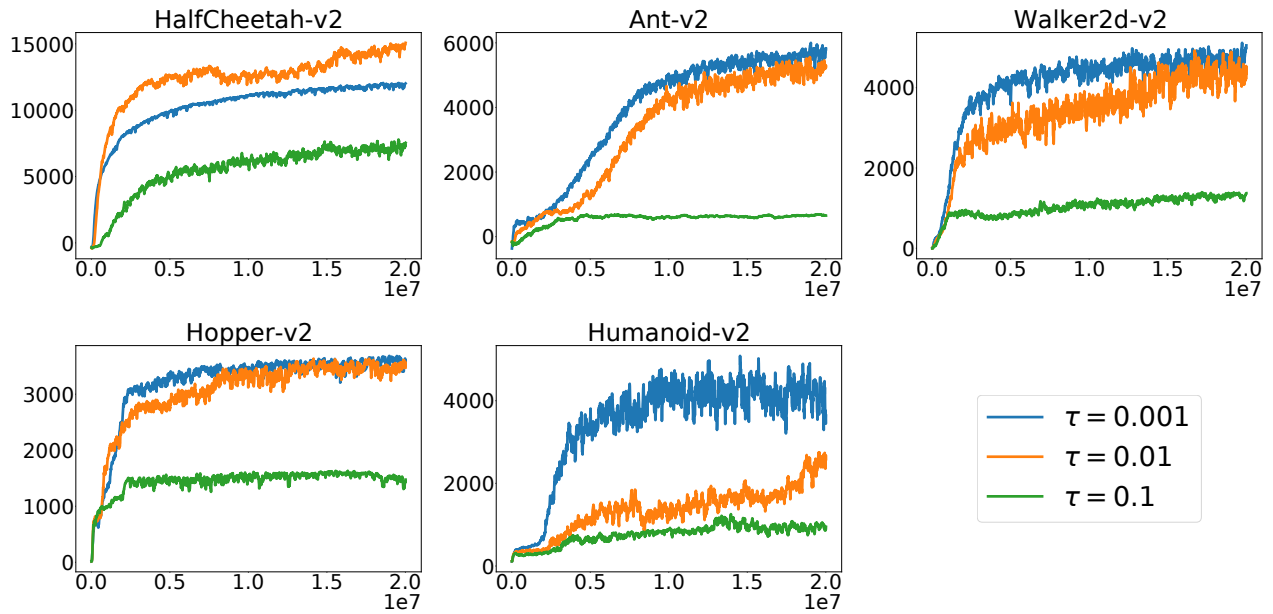


Figure 10: IQ ( $\alpha = 0$ ) with several values for  $\tau$  on all Gym tasks. Each line is the average over 10 seeds.

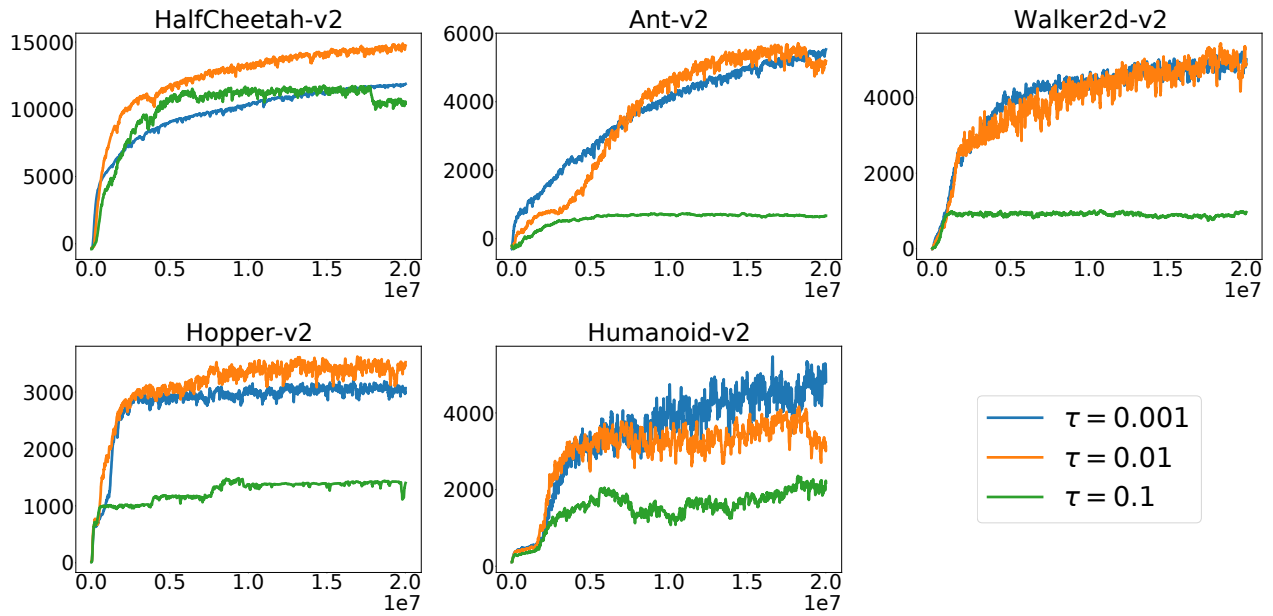


Figure 11: M-IQ ( $\alpha = 0.9$ ) with several values for  $\tau$  on all Gym tasks. Each line is the average over 10 seeds.

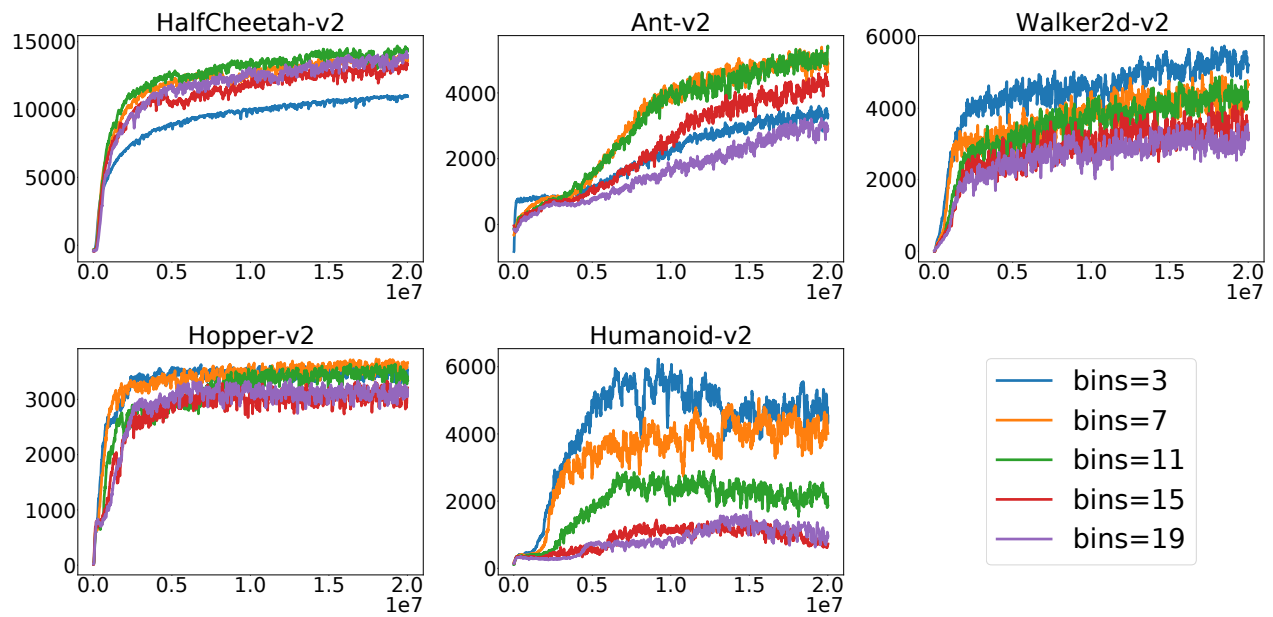


Figure 12: IQ with several number of bins on all Gym environments. Each line is the average over 10 seeds.