

---

# Randomized Prior Functions for Deep Reinforcement Learning

---

**Ian Osband**  
DeepMind  
iosband@google.com

**John Aslanides**  
DeepMind  
jaslanides@google.com

**Albin Cassirer**  
DeepMind  
cassirer@google.com

## Abstract

Dealing with uncertainty is essential for efficient reinforcement learning. There is a growing literature on uncertainty estimation for deep learning from fixed datasets, but many of the most popular approaches are poorly-suited to sequential decision problems. Other methods, such as bootstrap sampling, have no mechanism for uncertainty that does not come from the observed data. We highlight why this can be a crucial shortcoming and propose a simple remedy through addition of a randomized untrainable ‘prior’ network to each ensemble member. We prove that this approach is efficient with linear representations, provide simple illustrations of its efficacy with nonlinear representations and show that this approach scales to large-scale problems far better than previous attempts.

## 1 Introduction

Deep learning methods have emerged as the state of the art approach for many challenging problems [30, 70]. This is due to the statistical flexibility and computational scalability of large and deep neural networks, which allows them to harness the information in large and rich datasets. Deep reinforcement learning combines deep learning with sequential decision making under uncertainty. Here an agent takes actions inside an environment in order to maximize some cumulative reward [63]. This combination of deep learning with reinforcement learning (RL) has proved remarkably successful [67, 42, 60].

At the same time, elementary decision theory shows that the only admissible decision rules are Bayesian [12, 71]. Colloquially, this means that any decision rule that is not Bayesian can be improved (or even exploited) by some Bayesian alternative [14]. Despite this fact, the majority of deep learning research has evolved outside of Bayesian (or even statistical) analysis [55, 32]. This disconnect extends to deep RL, where the majority of state of the art algorithms have no concept of uncertainty [42, 41] and can fail spectacularly even in simple problems where success requires its consideration [40, 45].

There is a long history of research in Bayesian neural networks that never quite became mainstream practice [37, 43]. Recently, Bayesian deep learning has experienced a resurgence of interest with a myriad of approaches for uncertainty quantification in fixed datasets and also sequential decision problems [29, 11, 20, 47]. In this paper we highlight the surprising fact that many of these well-cited and popular methods for uncertainty estimation in deep learning can be poor choices for sequential decision problems. We show that this disconnect is more than a technical detail, but a serious shortcoming that can lead to arbitrarily poor performance. We support our claims by a series of simple lemmas for simple environments, together with experimental evidence in more complex settings.

Our approach builds on an alternative method for uncertainty in deep RL inspired by the statistical bootstrap [15]. This approach trains an ensemble of models, each on perturbed versions of the data. The resulting distribution *of the ensemble* is used to approximate the uncertainty in the estimate [47]. Although typically regarded as a frequentist method, bootstrapping gives near-optimal convergence rates when used as an approximate Bayesian posterior [19, 18]. However, these ensemble-based approaches to uncertainty quantification approximate a ‘posterior’ without an effective methodology to inject a ‘prior’. This can be a crucial shortcoming in sequential decision problems.

In this paper, we present a simple modification where each member of the ensemble is initialized together with a random but fixed *prior function*. Predictions in each ensemble member are then taken as the sum of the trainable neural network and its prior function. Learning/optimization is performed so that this sum (network plus prior) minimizes training loss. Therefore, with sufficient network capacity and optimization, the ensemble members will agree at observed data. However, in regions of the space with little or no training data, their predictions will be determined by the generalization of their networks and priors. Surprisingly, we show that this approach is equivalent to exact Bayesian inference for the special case of Gaussian linear models. Following on from this ‘sanity check’, we present a series of simple experiments designed to extend this intuition to deep learning. We show that many of the most popular approaches for uncertainty estimation in deep RL do *not* pass these sanity checks, and crystallize these shortcomings in a series of lemmas and small examples. We demonstrate that our simple modification can facilitate aspiration in difficult tasks where previous approaches for deep RL fail. We believe that this work presents a simple and practical approach to encoding prior knowledge with deep reinforcement learning.

## 2 Why do we need a ‘prior’ mechanism for deep RL?

We model the environment as a Markov decision process  $M = (\mathcal{S}, \mathcal{A}, R, P)$  [10]. Here  $\mathcal{S}$  is the state space and  $\mathcal{A}$  is the action space. At each time step  $t$ , the agent observes state  $s_t \in \mathcal{S}$ , takes action  $a_t \in \mathcal{A}$ , receives reward  $r_t \sim R(s_t, a_t)$  and transitions to  $s_{t+1} \sim P(s_t, a_t)$ . A policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  maps states to actions and let  $\mathcal{H}_t$  denote the history of observations before time  $t$ . An RL algorithm maps  $\mathcal{H}_t$  to a distribution over policies; we assess its quality through the cumulative reward over unknown environments. To perform well, an RL algorithm must learn to optimize its actions, combining both learning and control [63]. A ‘deep’ RL algorithm uses neural networks for nonlinear function approximation [32, 42].

The scale and scope of problems that might be approached through deep RL is vast, but there are three key aspects an efficient (and general) agent must address [63]:

1. **Generalization:** be able to learn from data it collects.
2. **Exploration:** prioritize the best experiences to learn from.
3. **Long-term consequences:** consider external effects beyond a single time step.

In this paper we focus on the importance of some form of ‘prior’ mechanism for efficient exploration. As a motivating example we consider a sparse reward task where random actions are very unlikely to ever see a reward. If an agent has never seen a reward then it is essential that some other form of aspiration, motivation, drive or curiosity direct its learning. We call this type of drive a ‘prior’ effect, since it does not come from the observed data, but are ambivalent as to whether this effect is philosophically ‘Bayesian’. Agents that do not have this prior drive will be left floundering aimlessly and thus may require exponentially large amounts of data in order to learn even simple problems [27].

To solve a specific task, it can be possible to attain superhuman performance without significant prior mechanism [42, 41]. However, if our goal is artificial *general* intelligence, then it is disconcerting that our best agents can perform very poorly even in simple problems [33, 39]. One potentially general approach to decision making is given by the Thompson sampling heuristic<sup>1</sup>: ‘randomly take action according to the probability you believe it is the optimal action’ [68]. In recent years there have been several attempts to apply this heuristic

---

<sup>1</sup>This heuristic is general in the sense that Thompson sampling can be theoretically justified in many of the domains where these other approaches fail [1, 48, 34, 58].

to deep reinforcement learning, each attaining significant outperformance over deep RL baselines on certain tasks [20, 47, 35, 11, 17]. In this section we outline crucial shortcomings for the most popular existing approaches to posterior approximation; these outlines will be brief, but more detail can be found in Appendix C. These shortcomings set the scene for Section 3, where we introduce a simple and practical alternative that passes each of our simple sanity checks: bootstrapped ensembles with randomized prior functions. In Section 4 we demonstrate that this approach scales gracefully to complex domains with deep RL.

## 2.1 Dropout as posterior approximation

One of the most popular modern approaches to regularization in deep learning is dropout sampling [61]. During training, dropout applies an independent random Bernoulli mask to the activations and thus guards against excessive co-adaptation of weights. Recent work has sought to understand dropout through a Bayesian lens, highlighting the connection to variational inference and arguing that the resultant dropout distribution approximates a Bayesian posterior [20]. This narrative has proved popular despite the fact that dropout distribution can be a poor approximation to most reasonable Bayesian posteriors [22, 46]:

**Lemma 1** (Dropout distribution does not concentrate with observed data).

*Consider any loss function  $\mathcal{L}$ , regularizer  $\mathcal{R}$  and data  $\mathcal{D}=\{(x,y)\}$ . Let  $\theta$  be parameters of any neural network architecture  $f$  trained with dropout rate  $p\in(0,1)$  and dropout masks  $W$ ,*

$$\theta_p^* \in \arg \min_{\theta} \mathbb{E}_{W \sim \text{Ber}(p), (x,y) \sim \mathcal{D}} [\mathcal{L}(x, y \mid \theta, W) + \mathcal{R}(\theta)]. \quad (1)$$

*Then the dropout distribution  $f_{\theta_p^*, W}$  is invariant to duplicates of the dataset  $\mathcal{D}$ .*

Lemma 1 is somewhat contrived, but highlights a clear shortcoming of dropout as posterior sampling: the dropout rate does not depend on the data. Lemma 1 means no agent employing dropout for posterior approximation can tell the difference between observing a set of data once and observing it  $N \gg 1$  times. This can lead to arbitrarily poor decision making, even when combined with an efficient strategy for exploration [45].

## 2.2 Variational inference and Bellman error

Dropout as posterior is motivated by its connection to variational inference (VI) [20], and recent work to address Lemma 1 improves the quality of this variational approximation by tuning the dropout rate from data [21].<sup>2</sup> However, there is a deeper problem to this line of research that is common across many works in this field: even given access to an oracle method for *exact* inference, applying independent inference to the Bellman error does not propagate uncertainty correctly for the value function as a whole [44]. To estimate the uncertainty in  $Q$  from the Bellman equation  $Q(s_t, a_t) = \mathbb{E}[r_{t+1} + \gamma \max_{\alpha} Q(s_{t+1}, \alpha)]$  it is crucial that the two sides of this equation are not independent random variables. Ignoring this dependence can lead to very bad estimates, even with exact inference.

**Lemma 2** (Independent VI on Bellman error does not propagate uncertainty).

*Let  $Y \sim N(\mu_Y, \sigma_Y^2)$  be a target value. If we train  $X \sim N(\mu, \sigma^2)$  according to the squared error*

$$\mu^*, \sigma^* \in \arg \min_{\mu, \sigma} \mathbb{E} [(X - Y)^2] \quad \text{for } X, Y \text{ independent}, \quad (2)$$

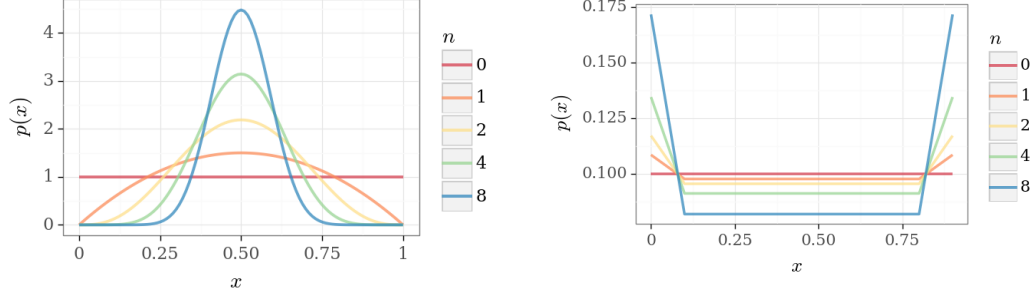
*then the solution  $\mu^* = \mu_Y, \sigma^* = 0$  propagates zero uncertainty from  $Y$  to  $X$ .*

To understand the significance of Lemma 2, imagine a deterministic system that transitions from  $s_1$  to  $s_2$  without reward. Suppose an agent is able to correctly quantify their posterior uncertainty for the value  $V(s_2) = Y \sim N(\mu_Y, \sigma_Y^2)$ . Training  $V(s_1) = X$  according to (2) will lead to zero uncertainty estimates at  $s_1$ , when in fact  $V(s_1) \sim N(\mu_Y, \sigma_Y^2)$ . This observation may appear simplistic, and may not say much more than ‘do not use the squared loss’ for VI in this setting. However, despite this clear failing (2) is precisely the loss used by the majority of approaches to VI for RL [17, 35, 65, 69, 20]. Note that this failure occurs even without decision making, function approximation and even when the true posterior lies within the variational class.

<sup>2</sup>Concrete dropout asymptotically improves the quality of the variational approximation, but provides no guarantees on its rate of convergence or error relative to exact Bayesian inference [21].

### 2.3 ‘Distributional reinforcement learning’

The key ingredient for a Bayesian formulation for sequential decision making is to consider beliefs not simply as a point estimate, but as a *distribution*. Recently an approach called ‘distributional RL’ has shown great success in improving stability and performance in deep RL benchmark algorithms [8]. Despite the name, these two ideas are quite distinct. ‘Distributional RL’ replaces a scalar estimate for the value function by a distribution that is trained to minimize a loss against the distribution of data it observes. This distribution of observed data is an orthogonal concept to that of Bayesian uncertainty.



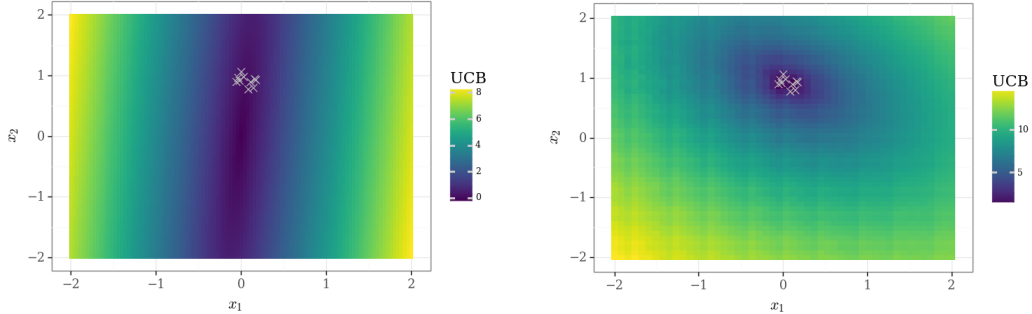
(a) Posterior beliefs concentrate around  $p = 0.5$ . (b) ‘Distributional’ tends to mass at 0 and 1.  
Figure 1: Output distribution after observing  $n$  heads and  $n$  tails of a coin.

Figure 1 presents an illustration of these two distributions after observing flips of a coin. As more data is gathered the posterior distribution concentrates around the mean whereas the ‘distributional RL’ estimate approaches that of the generating Bernoulli. Although both approaches might reasonably claim a ‘distributional perspective’ on RL, these two distributions have orthogonal natures and behave quite differently. Conflating one for the other can lead to arbitrarily poor decisions; it is the uncertainty in beliefs (‘epistemic’), not the distributional noise (‘aleatoric’) that is important for exploration [27].

### 2.4 Count-based uncertainty estimates

Another popular method for incentivizing exploration is with a density model or ‘pseudocount’ [6]. Inspired by the analysis of tabular systems, these models assign a bonus to states and actions that have been visited infrequently according to a density model. This method can perform well, but only when the generalization of the density model is aligned with the task objective. Crucially, this generalization is not learned from the task [53].

Even with an optimal state representation and density, a count-based bonus on states can be poorly suited for efficient exploration. Consider a linear bandit with reward  $r_t(x_t) = x_t^T \theta^* + \epsilon_t$  for some  $\theta^* \in \mathbb{R}^d$  and  $\epsilon_t \sim N(0, 1)$  [56]. Figure 2 compares the uncertainty in the expected reward  $\mathbb{E}[x^T \theta^*]$  with that obtained by density estimation on the observed  $x_t$ . A bonus based upon the state density does not correlate with the *uncertainty* over the unknown optimal action. This disconnect can lead to arbitrarily poor decisions [49].



(a) Uncertainty bonus from posterior over  $x^T \theta^*$ . (b) Bonus from Gaussian pseudocount  $p(x)$ .  
Figure 2: Count-based uncertainty leads to a poorly aligned bonus even in a linear system.

### 3 Randomized prior functions for deep ensembles

Section 2 motivates the need for effective uncertainty estimates in deep RL. We note that crucial failure cases of several popular approaches can arise even with simple linear models. As a result, we take a moment to review the setting of Bayesian linear regression. Let  $\theta \in \mathbb{R}^d$  with prior  $N(\bar{\theta}, \lambda I)$  and data  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  for  $x_i \in \mathbb{R}^d$  and  $y_i = \theta^T x_i + \epsilon_i$  with  $\epsilon_i \sim N(0, \sigma^2)$  iid. Then, conditioned on  $\mathcal{D}$ , the posterior for  $\theta$  is Gaussian:

$$\mathbb{E}[\theta|\mathcal{D}] = \left( \frac{1}{\sigma^2} X^T X + \frac{1}{\lambda} I \right)^{-1} \left( \frac{1}{\sigma^2} X^T y + \frac{1}{\lambda} \bar{\theta} \right), \quad \text{Cov}[\theta|\mathcal{D}] = \left( \frac{1}{\sigma^2} X^T X + \frac{1}{\lambda} I \right)^{-1}. \quad (3)$$

Equation (3) relies on Gaussian conjugacy and linear models, which cannot easily be extended to deep neural networks. The following result shows that we can replace this analytical result with a simple computational procedure.

**Lemma 3** (Computational generation of posterior samples).

Let  $f_\theta(x) = x^T \theta$ ,  $\tilde{y}_i \sim N(y_i, \sigma^2)$  and  $\tilde{\theta} \sim N(\bar{\theta}, \lambda I)$ . Then the either of the following optimization problems generate a sample  $\theta | \mathcal{D}$  according to (3):

$$\arg \min_{\theta} \sum_{i=1}^n \|\tilde{y}_i - f_\theta(x_i)\|_2^2 + \frac{\sigma^2}{\lambda} \|\tilde{\theta} - \theta\|_2^2, \quad (4)$$

$$\tilde{\theta} + \arg \min_{\theta} \sum_{i=1}^n \|\tilde{y}_i - (f_{\tilde{\theta}} + f_\theta)(x_i)\|_2^2 + \frac{\sigma^2}{\lambda} \|\theta\|_2^2. \quad (5)$$

*Proof.* To prove (4) note that the solution is Gaussian and then match moments; equation (5) then follows by relabeling [49].  $\square$

Lemma 3 is revealing since it allows us to view Bayesian regression through a purely computational lens: ‘generate posterior samples by training on noisy versions of the data, together with some random regularization’. Even for nonlinear  $f_\theta$ , we can still compute (4) or (5). Although the resultant  $f_\theta$  will no longer be an exact posterior, at least it passes the ‘sanity check’ in this simple linear setting (unlike the approaches of Section 2). We argue this method is quite intuitive: the perturbed data  $\tilde{\mathcal{D}} = \{(x_i, \tilde{y}_i)\}_{i=1}^n$  is generated according to the estimated noise process  $\epsilon_t$  and the sample  $\tilde{\theta}$  is drawn from prior beliefs. Intuitively (4) says to fit to  $\tilde{\mathcal{D}}$  and regularize weights to a prior sample of weights  $\tilde{\theta}$ ; (5) says to generate a prior function  $f_{\tilde{\theta}}$  and then fit an additive term to noisy data  $\tilde{\mathcal{D}}$  with regularized complexity.

This paper explores the performance of each of these methods for uncertainty estimation with deep learning. We find empirical support that method (5) coupled with a *randomized prior function* significantly outperforms ensemble-based approaches without prior mechanism. We also find that (5) significantly outperforms (4) in deep RL. We suggest a major factor in this comes down to the huge dimensionality of neural network weights, whereas the output policy or value is typically far smaller. In this case, it makes sense to enforce prior beliefs in the low dimensional space. Further, the initialization of neural network weights plays an important role in their generalization properties and optimization via stochastic gradient descent (SGD) [23, 38]. As such, (5) may help to decouple the dual roles of initial weights as both ‘prior’ and training initializer. Algorithm 1 describes our approach applied to modern deep learning architectures.

---

#### Algorithm 1 Randomized prior functions for ensemble posterior.

---

**Require:** Data  $\mathcal{D} \subseteq \{(x, y) | x \in \mathcal{X}, y \in \mathcal{Y}\}$ , loss function  $\mathcal{L}$ , neural model  $f_\theta: \mathcal{X} \rightarrow \mathcal{Y}$ , Ensemble size  $K \in \mathbb{N}$ , noise procedure `data_noise`, distribution over priors  $\mathcal{P} \subseteq \{\mathbb{P}(p) | p: \mathcal{X} \rightarrow \mathcal{Y}\}$ .

- 1: **for**  $k = 1, \dots, K$  **do**
- 2:     initialize  $\theta_k \sim$  Glorot initialization [23].
- 3:     form  $\mathcal{D}_k = \text{data\_noise}(\mathcal{D})$  (e.g. Gaussian noise or bootstrap sampling [50]).
- 4:     sample prior function  $p_k \sim \mathcal{P}$ .
- 5:     optimize  $\nabla_{\theta | \theta = \theta_k} \mathcal{L}(f_\theta + p_k; \mathcal{D}_k)$  via ADAM [28].
- 6: **return** ensemble  $\{f_{\theta_k} + p_k\}_{k=1}^K$ .

---

## 4 Deep reinforcement learning

Algorithm 1 might be applied to model or policy learning approaches, but this paper focuses on value learning. We apply Algorithm 1 to *deep Q networks* (DQN) [42] on a series of tasks designed to require good uncertainty estimates. We train an ensemble of  $K$  networks  $\{Q_k\}_{k=1}^K$  in parallel, each on a perturbed version of the observed data  $\mathcal{H}_t$  and each with a distinct random, but fixed, prior function  $p_k$ . Each episode, the agent selects  $j \sim \text{Unif}([1, \dots, K])$  and follows the greedy policy w.r.t.  $Q_j$  for the duration of the episode. This algorithm is essentially bootstrapped DQN (BootDQN) except for the addition of the prior function  $p_k$  [47]. We use the statistical bootstrap rather than Gaussian noise (5) to implement a state-specific variance [19].

Let  $\gamma \in [0, 1]$  be a discount factor that induces a time preference over future rewards. For a neural network family  $f_\theta$ , prior function  $p$ , and data  $\mathcal{D} = \{(s_t, a_t, r_t, s'_t)\}$  we define the  $\gamma$ -discounted empirical temporal difference (TD) loss,

$$\mathcal{L}_\gamma(\theta; \theta^-, p, \mathcal{D}) := \sum_{t \in \mathcal{D}} \left( r_t + \gamma \overbrace{\max_{a' \in \mathcal{A}} (f_{\theta^-} + p)(s'_t, a')}^{\text{target } Q} - \overbrace{(f_\theta + p)(s_t, a_t)}^{\text{online } Q} \right)^2. \quad (6)$$

Using this notation, the learning update for BootDQN with prior functions is a simple application of Algorithm 1, which we outline below. To complete the RL algorithm we implement a 50-50 `ensemble_buffer`, where each transition has a 50% chance of being included in the replay for model  $k = 1, \dots, K$ . For a complete description of BootDQN+prior agent, see Appendix A.

---

### Algorithm 2 learn\_bootstrapped\_dqn\_with\_prior

---

<b>Agent:</b>	$\theta_1, \dots, \theta_K$	trainable network parameters
	$p_1, \dots, p_K$	fixed prior functions
	$\mathcal{L}_\gamma(\theta = \cdot; \theta^- = \cdot, p = \cdot, \mathcal{D} = \cdot)$	TD error loss function
	<code>ensemble_buffer</code>	replay buffer of $K$ -parallel perturbed data
<b>Updates:</b>	$\theta_1, \dots, \theta_K$	agent value function estimate
1:	<b>for</b> $k$ in $(1, \dots, K)$ <b>do</b>	
2:	Data $\mathcal{D}_k \leftarrow \text{ensemble\_buffer}[k].\text{sample\_minibatch}()$	
3:	optimize $\nabla_{\theta \theta=\theta_k} \mathcal{L}(\theta; \theta_k, p_k, \mathcal{D}_k)$ via ADAM [28].	

---

#### 4.1 Does BootDQN+prior address the shortcomings from Section 2?

Algorithm 2 is a simple modification of vanilla Q-learning: rather than maintain a single point estimate for  $Q$ , we maintain  $K$  estimates in parallel, and rather than regularize each estimate to a single value, each is individually regularized to a distinct random prior function. We show that that this simple and scalable algorithm overcomes the crucial shortcomings that afflict existing methods, as outlined in Section 2.

- ✓ **Posterior concentration** (Section 2.1): Prior function + noisy data means the ensemble is initially diverse, but concentrates as more data is gathered. For linear-gaussian systems this matches Bayes posterior, bootstrap offers a general, non-parametric approach [16, 18].
- ✓ **Multi-step uncertainty** (Section 2.2): Since each network  $k$  trains only on its *own* target value, BootDQN+prior propagates a temporally-consistent sample of  $Q$ -value [49].
- ✓ **Epistemic vs aleatoric** (Section 2.3): BootDQN+prior optimises the *mean* TD loss (6) and does not seek to fit the noise in returns, unlike ‘distributional RL’ [7].
- ✓ **Task-appropriate generalization** (Section 2.4): We explore according to our uncertainty in the value  $Q$ , rather than density on state. As such, our generalization naturally occurs in the space of *features* relevant to the task, rather than pixels or noise [6].
- ✓ **Intrinsic motivation** (comparison to BootDQN without prior): In an environment with zero rewards, a bootstrap ensemble may simply learn to predict zero for *all* states. The prior  $p_k$  can make this generalization unlikely for  $Q_k$  at unseen states  $\tilde{s}$  so  $\mathbb{E}[\max_\alpha Q_k(\tilde{s}, \alpha)] > 0$ ; thus BootDQN+prior seeks novelty even with no observed rewards.

Another source of justification comes from the observation that BootDQN+prior is an instance of *randomized least-squares value iteration* (RLSVI), with regularization via ‘prior

function’ for an ensemble of neural networks. RLSVI with linear function approximation and Gaussian noise guarantees a bound on expected regret of  $\tilde{O}(\sqrt{|S||A|T})$  in the tabular setting [49].<sup>3</sup> Similarly, analysis for the bandit setting establishes that  $K = \tilde{O}(|A|)$  models trained online can attain similar performance to full resampling each episode [36]. Our work in this paper pushes beyond the boundaries of these analyses, which are presented as ‘sanity checks’ that our algorithm is at least sensible in simple settings, rather than a certificate of correctness for more complex ones. The rest of this paper is dedicated to an empirical investigation of our algorithm through computational experiments. Encouragingly, we find that many of the insights born out of simple problems extend to more complex ‘deep RL’ settings and good evidence for the efficacy of our algorithm.

## 4.2 Computational experiments

Our experiments focus on a series of environments that require deep exploration together with increasing state complexity [27, 49]. In each of our domains, random actions are very unlikely to achieve a reward and exploratory actions may even come at a cost. Any algorithm without prior motivation will have no option but to explore randomly, or worse, eschew exploratory actions completely in the name of premature and sub-optimal exploitation. In our experiments we focus on a *tabula rasa* setting in which the prior function is drawn as a random neural network. Although our prior distribution  $\mathcal{P}$  could encode task-specific knowledge (e.g. through sampling the true  $Q^*$ ), we leave this investigation to future work.

### 4.2.1 Chain environments

We begin our experiments with a family of chain-like environments that highlight the need for deep exploration [62]. The environments are indexed by problem scale  $N \in \mathbb{N}$  and action mask  $W \sim \text{Ber}(0.5)^{N \times N}$ , with  $\mathcal{S} = \{0,1\}^{N \times N}$  and  $\mathcal{A} = \{0,1\}$ . The agent begins each episode in the upper left-most state in the grid and deterministically falls one row per time step. The state encodes the agent’s row and column as a one-hot vector  $s_t \in \mathcal{S}$ . The actions  $\{0,1\}$  move the agent left or right depending on the action mask  $W$  at state  $s_t$ , which remains fixed. The agent incurs a cost of  $0.01/N$  for moving right in all states except for the right-most, in which the reward is 1. The reward for action left is always zero. An episode ends after  $N$  time steps so that the optimal policy is to move right each step and receive a total return of 0.99; all other policies receive zero or negative return. Crucially, algorithms without deep exploration take  $\Omega(2^N)$  episodes to learn the optimal policy [52].<sup>4</sup>

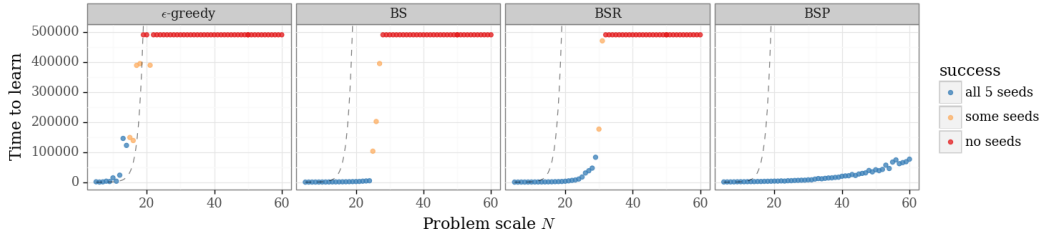


Figure 3: Only bootstrap with additive prior network (BSP) scales gracefully to large problems. Plotting BSP on a log-log scale suggests an empirical scaling  $T_{\text{learn}} = \tilde{O}(N^3)$ ; see Figure 8.

Figure 3 presents the average time to learn for  $N = 5, \dots, 60$  up to 500K episodes over 5 seeds and ensemble  $K = 20$ . We say that an agent has learned the optimal policy when the average regret per episode drops below 0.9. We compare three variants of BootDQN, depending on their mechanism for ‘prior’ effects. **BS** is bootstrap without prior mechanism. **BSR** is bootstrap with  $l_2$  regularization on weights per (4). **BSP** is bootstrap with additive prior function per (5). In each case we initialize a random 20-unit MLP; BSR regularizes to these initial weights and BSP trains an additive network. Although all bootstrap methods significantly outperform  $\epsilon$ -greedy, only BSP successfully scales to large problem sizes.

Figure 4 presents a more detailed analysis of the sensitivity of our approach to the tuning parameters of different regularization approaches. We repeat the experiments of Figure 3

<sup>3</sup>Regret measures the shortfall in cumulative rewards compared to that of the optimal policy.

<sup>4</sup>The dashed lines indicate the  $2^N$  dithering lower bound. The action mask  $W$  means this cannot be solved easily by evolution or policy search evolution, unlike previous ‘chain’ examples [47, 54].

and examine the size of the largest problem solved before 50K episodes. In each case larger ensembles lead to better performance, but this effect tends to plateau relatively early. Figure 4a shows that regularization provides little or no benefit to BSR. Figure 4b examines the effect of scaling the randomly initialized MLP by a scalar hyperparameter  $\beta$ .

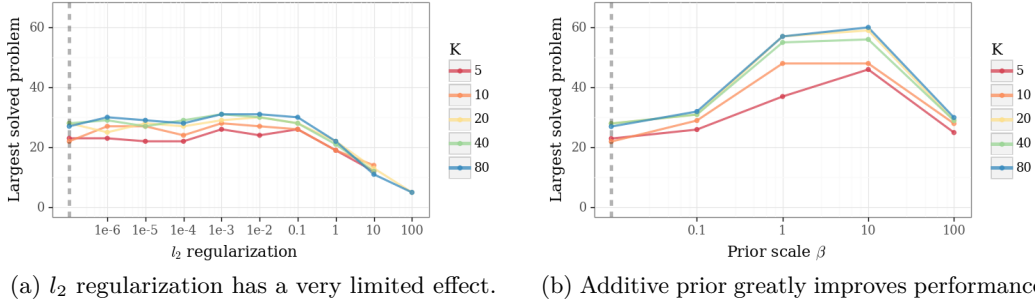


Figure 4: Comparing effects of different styles of prior regularization in Bootstrapped DQN.

### How does BSP solve this exponentially-difficult problem?

At first glance this ‘chain’ problem may seem like an impossible task. Finding the single rewarding policy out of  $2^N$  is not simply a needle-in-a-haystack, but more akin to looking for a piece of hay in a needle-stack! Since every policy apart from the rewarding one is painful, it’s very tempting for an agent to give up and receive reward zero. We now provide some intuition for how BSP is able to consistently and scalably solve such a difficult task.

One way to interpret this result is through analysing BSP with linear function approximation via Lemma 3. As outlined in Section 4.1, BSP with linear function approximation satisfies a polynomial regret bound [49]. Further, this empirical scaling matches that predicted by the regret bound tabular domain [51] (see Figure 8). Here, the prior function plays a crucial role - it provides motivation for the agent to explore even when the observed data has low (or no) reward. Note that it is not necessary the sampled prior function leads to a good policy itself; in fact this is exponentially unlikely according to our initialization scheme. The crucial element is that when a new state  $s'$  is reached there is *some* ensemble member that estimates  $\max_{a'} Q_k(s', a')$  is sufficiently positive to warrant visiting, even if it causes some negative reward along the way. In that case, when network  $k$  is active it will seek out the potentially-informative  $s'$  even if it is multiple timesteps away; this effect is sometimes called *deep exploration*. We present an accompanying visualization at [http://bit.ly/rpf\\_nips](http://bit.ly/rpf_nips).

However, this connection to linear RLSVI does not inform why BSP should outperform BSR. To account for this, we appeal to the functional dynamics of deep learning architectures (see Section 3). In large networks weight decay (per BSR) may be an ineffective mechanism on the *output*  $Q$ -values. Instead, training an additive network via SGD (per BSP) may provide a more effective regularization on the output function [73, 38, 5]. We expand on this hypothesis and further details of these experiments in Appendix B.1. This includes investigation of NoisyNets [17] and dropout [20], which both perform poorly, and a comparison to UCB-based algorithms, which scale much worse than BSP, even with oracle access to state visit counts.

#### 4.2.2 Cartpole swing-up

The experiments of Section 4.2.1 show that the choice of prior mechanism can be absolutely essential for efficient exploration via randomized value functions. However, since the underlying system is a small finite MDP we might observe similar performance through a tabular algorithm. In this section we investigate a classic benchmark problem that necessitates nonlinear function approximation: cartpole [63]. We modify the classic formulation so that the pole begins hanging down and the agent only receives a reward when the pole is upright, balanced, and centered<sup>5</sup>. We also add a cost of 0.1 for moving the cart. This problem embodies many of the same aspects of 4.2.1, but since the agent interacts with the environment through state  $s_t = (\cos(\theta_t), \sin(\theta_t), \theta_t, x_t, \dot{x}_t)$ , the agent must also learn nonlinear generalization. Tabular approaches are not practical due to the curse of dimensionality.

<sup>5</sup>We use the DeepMind control suite [66] with reward +1 only when  $\cos(\theta) > 0.95$ ,  $|x| < 0.1$ ,  $|\dot{\theta}| < 1$ , and  $|\dot{x}| < 1$ . Each episode lasts 1,000 time steps, simulating 10 seconds of interaction.



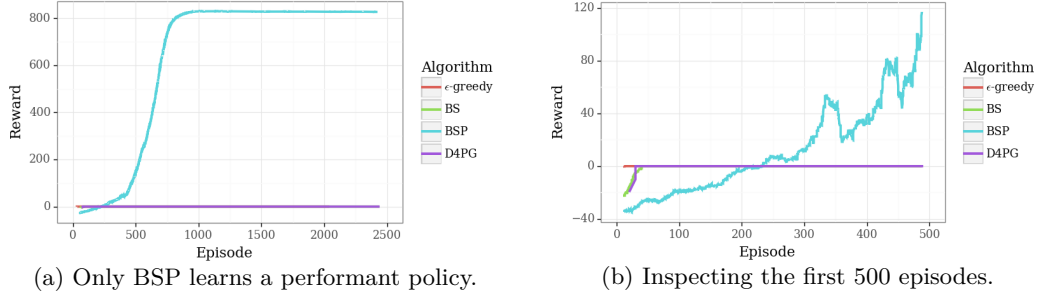


Figure 5: Learning curves for the modified cartpole swing-up task.

Figure 5 compares the performance of DQN with  $\epsilon$ -greedy, bootstrap without prior (BS), bootstrap with prior networks (BSP) and the state-of-the-art continuous control algorithm D4PG, itself an application of ‘distributional RL’ [4]. Only BSP learns a performant policy; no other approach ever attains any positive reward. We push experimental details, including hyperparameter analysis, to Appendix B.2. These results are significant in that they show that our intuitions translate from simple domains to more complex nonlinear settings, although the underlying state is relatively low dimensional. Our next experiments investigate performance in a high dimensional and visually rich domain.

#### 4.2.3 Montezuma’s revenge

Our final experiment comes from the Arcade Learning Environment and the canonical sparse reward game, Montezuma’s Revenge [9]. The agent interacts directly with the pixel values and, even under an optimal policy, there can be hundreds of time steps between rewarding actions. This problem presents a significant exploration challenge in a visually rich environment; many published algorithms are essentially unable to attain any reward here [42, 41]. We compare performance against a baseline distributed DQN agent with double Q-learning, prioritized experience replay and dueling networks [25, 24, 59, 72]. To save computation we follow previous work and use a shared convnet for the ensemble uncertainty [47, 3]. Figure 6 presents the results for varying prior scale  $\beta$  averaged over three seeds. Once again, we see that the prior network can be absolutely critical to successful exploration.

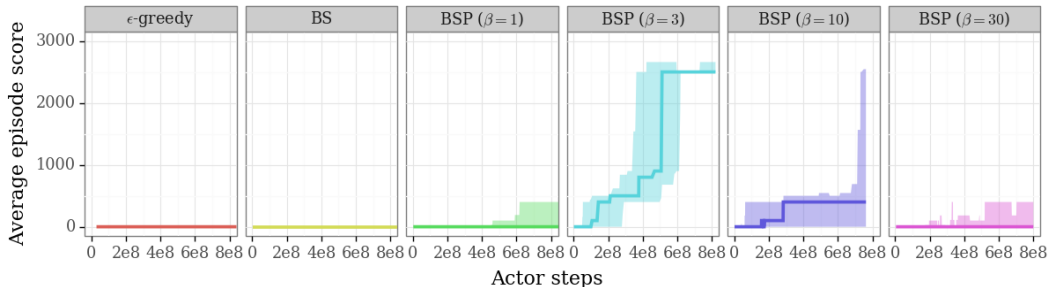


Figure 6: The prior network qualitatively changes behavior on Montezuma’s revenge.

## 5 Conclusion

This paper highlights the importance of uncertainty estimates in deep RL, the need for an effective ‘prior’ mechanism, and its potential benefits towards efficient exploration. We present some alarming shortcomings of existing methods and suggest bootstrapped ensembles with randomized prior functions as a simple, practical alternative. We support our claims through an analysis of this method in the linear setting, together with a series of simple experiments designed to highlight the key issues. Our work leaves several open questions. What kinds of prior functions are appropriate for deep RL? Can they be optimized or ‘meta-learned’? Can we distill the ensemble process to a single network? We hope this work helps to inspire solutions to these problems, and also build connections between the theory of efficient learning and practical algorithms for deep reinforcement learning.

## Acknowledgements

We would like to thank many people who made important contributions to this paper. This paper can be thought of as a specific type of ‘deep exploration via randomized value functions’, whose line of research has been crucially driven by the contributions of (and conversations with) Benjamin Van Roy, Daniel Russo and Zheng Wen. Further, we would like to acknowledge the many helpful comments and support from Mohammad Gheshlaghi Azar, David Budden, David Silver and Justin Sirignano. Finally, we would like to make a special mention for Hado Van Hasselt, who coined the term ‘hay in a needle-stack’ to describe our experiments from Section 4.2.1.

## References

- [1] Shipra Agrawal and Navin Goyal. Analysis of Thompson sampling for the multi-armed bandit problem. In *Conference on Learning Theory*, pages 39–1, 2012.
- [2] Shipra Agrawal and Navin Goyal. Further optimal regret bounds for Thompson sampling. In *Artificial Intelligence and Statistics*, pages 99–107, 2013.
- [3] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration through bayesian deep q-networks. *arXiv preprint arXiv:1802.04412*, 2018.
- [4] Gabriel Barth-Maron, Matthew W Hoffman, David Budden, Will Dabney, Dan Horgan, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional deterministic policy gradients. *arXiv preprint arXiv:1804.08617*, 2018.
- [5] Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. In *Advances in Neural Information Processing Systems 30*, pages 6241–6250, 2017.
- [6] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Rémi Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems 29*, pages 1471–1479, 2016.
- [7] Marc G Bellemare, Will Dabney, and Rémi Munos. A Distributional Perspective on Reinforcement Learning. *Proceedings of the 34th International Conference on Machine Learning (ICML)*, 2017.
- [8] Marc G Bellemare, Will Dabney, and Rémi Munos. A distributional perspective on reinforcement learning. In *International Conference on Machine Learning*, pages 449–458, 2017.
- [9] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. The arcade learning environment: An evaluation platform for general agents. *J. Artif. Intell. Res. (JAIR)*, 47:253–279, 2013.
- [10] Dimitri P. Bertsekas and John Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, September 1996.
- [11] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- [12] David Roxbee Cox and David Victor Hinkley. *Theoretical statistics*. CRC Press, 1979.
- [13] Will Dabney, Mark Rowland, Marc G Bellemare, and Rémi Munos. Distributional reinforcement learning with quantile regression. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [14] Bruno De Finetti. La prévision: ses lois logiques, ses sources subjectives. In *Annales de l’institut Henri Poincaré*, volume 7, pages 1–68, 1937.
- [15] Bradley Efron. *The jackknife, the bootstrap and other resampling plans*, volume 38. SIAM, 1982.
- [16] Bradley Efron and Robert J Tibshirani. *An introduction to the bootstrap*. CRC press, 1994.
- [17] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Rémi Munos, Demis Hassabis, Olivier Pietquin, et al. Noisy networks for exploration. In *Proc. of ICLR*, 2018.
- [18] Tadayoshi Fushiki. Bootstrap prediction and bayesian prediction under misspecified models. *Bernoulli*, pages 747–758, 2005.
- [19] Tadayoshi Fushiki, Fumiyasu Komaki, Kazuyuki Aihara, et al. Nonparametric bootstrap prediction. *Bernoulli*, 11(2):293–307, 2005.
- [20] Yarin Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep learning. In *International Conference on Machine Learning*, 2016.

- [21] Yarın Gal, Jiri Hron, and Alex Kendall. Concrete dropout. In *Advances in Neural Information Processing Systems*, pages 3584–3593, 2017.
- [22] Yarın Gal, Rowan McAllister, and Carl Edward Rasmussen. Improving pilco with bayesian neural network dynamics models. In *Data-Efficient Machine Learning workshop, ICML*, 2016.
- [23] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the 13th international conference on artificial intelligence and statistics*, pages 249–256, 2010.
- [24] Hado van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, AAAI’16, pages 2094–2100. AAAI Press, 2016.
- [25] Daniel Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado Van Hasselt, and David Silver. Distributed prioritized experience replay. In *6th International Conference on Learning Representations*, 2018.
- [26] Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(Apr):1563–1600, 2010.
- [27] M. Kearns and S. Singh. Near-optimal reinforcement learning in polynomial time. *Machine Learning*, 49, 2002.
- [28] Diederik Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *Proceedings of the International Conference on Learning Representations*, 2015.
- [29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *International Conference on Learning Representations*, 2014.
- [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems 25*, pages 1097–1105, 2012.
- [31] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pages 6405–6416, 2017.
- [32] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.
- [33] Shane Legg, Marcus Hutter, et al. A collection of definitions of intelligence. *Frontiers in Artificial Intelligence and applications*, 157:17, 2007.
- [34] Jan Leike, Tor Lattimore, Laurent Orseau, and Marcus Hutter. Thompson sampling is asymptotically optimal in general environments. *Uncertainty in Artificial Intelligence*, 2016.
- [35] Zachary C Lipton, Jianfeng Gao, Lihong Li, Xiujun Li, Faisal Ahmed, and Li Deng. Efficient exploration for dialogue policy learning with bbq networks & replay buffer spiking. *arXiv preprint arXiv:1608.05081*, 2016.
- [36] Xiuyuan Lu and Benjamin Van Roy. Ensemble sampling. In *Advances in Neural Information Processing Systems*, pages 3260–3268, 2017.
- [37] David JC MacKay. A practical Bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- [38] Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient descent quantizes ReLU network features. *arXiv preprint arXiv:1803.08367*, 2018.
- [39] Horia Mania, Aurelia Guy, and Benjamin Recht. Simple random search provides a competitive approach to reinforcement learning. *arXiv preprint arXiv:1803.07055*, 2018.
- [40] Oliver Mihatsch and Ralph Neuneier. Risk-sensitive reinforcement learning. *Machine learning*, 49(2-3):267–290, 2002.
- [41] Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *Proc. of ICML*, 2016.
- [42] 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):529–533, 2015.
- [43] Radford M Neal. *Bayesian learning for neural networks*, volume 118. Springer Science & Business Media, 2012.
- [44] Brendan O’Donoghue, Ian Osband, Remi Munos, and Volodymyr Mnih. The uncertainty bellman equation and exploration. *arXiv preprint arXiv:1709.05380*, 2017.

- [45] Ian Osband. *Deep Exploration via Randomized Value Functions*. PhD thesis, Stanford University, 2016.
- [46] Ian Osband. Risk versus uncertainty in deep learning: Bayes, bootstrap and the dangers of dropout. 2016.
- [47] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via bootstrapped DQN. In *Advances In Neural Information Processing Systems 29*, pages 4026–4034, 2016.
- [48] Ian Osband, Daniel Russo, and Benjamin Van Roy. (More) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems 26*, pages 3003–3011. 2013.
- [49] Ian Osband, Daniel Russo, Zheng Wen, and Benjamin Van Roy. Deep exploration via randomized value functions. *arXiv preprint arXiv:1703.07608*, 2017.
- [50] Ian Osband and Benjamin Van Roy. Bootstrapped Thompson sampling and deep exploration. *arXiv preprint arXiv:1507.00300*, 2015.
- [51] Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for reinforcement learning? In *Proceedings of the 34th International Conference on Machine Learning*, pages 2701–2710, 2017.
- [52] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. In *Proceedings of The 33rd International Conference on Machine Learning*, pages 2377–2386, 2016.
- [53] Georg Ostrovski, Marc G Bellemare, Aaron van den Oord, and Rémi Munos. Count-based exploration with neural density models. In *Proc. of ICML*, 2017.
- [54] Matthias Plappert, Rein Houthoofd, Prafulla Dhariwal, Szymon Sidor, Richard Y Chen, Xi Chen, Tamim Asfour, Pieter Abbeel, and Marcin Andrychowicz. Parameter space noise for exploration. *arXiv preprint arXiv:1706.01905*, 2017.
- [55] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, DTIC Document, 1985.
- [56] Paat Rusmevichientong and John N. Tsitsiklis. Linearly parameterized bandits. *Math. Oper. Res.*, 35(2):395–411, 2010.
- [57] Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics of Operations Research*, 39(4):1221–1243, 2014.
- [58] Daniel Russo, Benjamin Van Roy, Abbas Kazerouni, and Ian Osband. A tutorial on Thompson sampling. *arXiv preprint arXiv:1707.02038*, 2017.
- [59] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *CoRR*, abs/1511.05952, 2015.
- [60] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- [61] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [62] Malcolm Strens. A Bayesian framework for reinforcement learning. In *International Conference on Machine Learning*, pages 943–950, 2000.
- [63] Richard Sutton and Andrew Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2017.
- [64] Richard S Sutton, Joseph Modayil, Michael Delp, Thomas Degris, Patrick M Pilarski, Adam White, and Doina Precup. Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *The 10th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*, pages 761–768. International Foundation for Autonomous Agents and Multiagent Systems, 2011.
- [65] Yunhao Tang and Alp Kucukelbir. Variational deep q network. *arXiv preprint arXiv:1711.11225*, 2017.
- [66] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, et al. Deepmind control suite. *arXiv preprint arXiv:1801.00690*, 2018.

- [67] Gerald Tesauro. Temporal difference learning and TD-gammon. *Communications of the ACM*, 38(3):58–68, 1995.
- [68] William R Thompson. On the likelihood that one unknown probability exceeds another in view of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- [69] Ahmed Touati, Harsh Satija, Joshua Romoff, Joelle Pineau, and Pascal Vincent. Randomized value functions via multiplicative normalizing flows. *arXiv preprint arXiv:1806.02315*, 2018.
- [70] Aaron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [71] Abraham Wald. Statistical decision functions. In *Breakthroughs in Statistics*, pages 342–357. Springer, 1992.
- [72] Ziyu Wang, Nando de Freitas, and Marc Lanctot. Dueling network architectures for deep reinforcement learning. *CoRR*, abs/1511.06581, 2015.
- [73] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. *CoRR*, abs/1611.03530, 2016.

## A Reinforcement learning algorithm

In this appendix we fill out the details for the complete pseudocode for the BootDQN+Prior RL agent. Our problem setting matches the description of prior work [49]; we reproduce the algorithms and figures in this section for the convenience of our readers. At a high level, our agents interact with the environment through repeated finite episodes as described by Algorithm 3. To describe an agent, we must simply implement the `act`, `update_buffer` and `learn_from_buffer` methods.

---

### Algorithm 3 live

---

**Input:**    `agent`                    methods `act`, `update_buffer`, `learn_from_buffer`  
              `environment`        methods `reset`, `step`

```

1: for episode = 1, 2, ... do
2:   agent.learn_from_buffer()
3:   transition  $\leftarrow$  environment.reset()
4:   while transition.new_state is not null do
5:     action  $\leftarrow$  agent.act(transition.new_state)
6:     transition  $\leftarrow$  environment.step(action)
7:     agent.update_buffer(transition)

```

---

BootDQN+prior implements an `ensemble_buffer` that maintains  $K$  buffers in parallel, although this may clearly be implemented in an efficient way that uses  $o(K)$  memory. Figure 7 provides an illustration of how BootDQN learns and maintains  $K$  estimates of the value function in parallel.

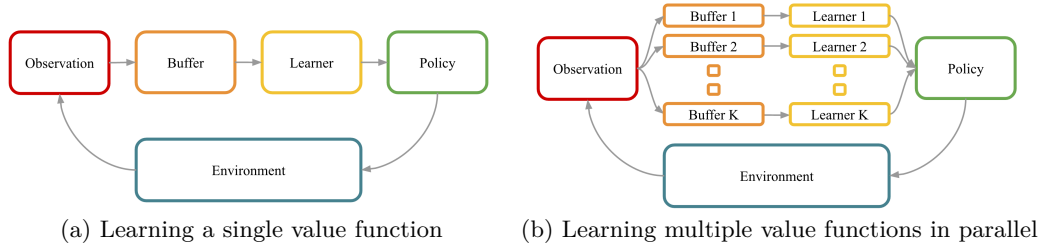


Figure 7: RLSVI via ensemble sampling, each member produced by LSVI on perturbed data.

To implement an online double-or-nothing bootstrap we employ Algorithm 4, which assigns each transition to each ensemble buffer with probability  $\frac{1}{2}$ .

---

### Algorithm 4 `ensemble_buffer.update_bootstrap(·)`

---

**Input:**        `transition`                     $(s_t, a_t, r_t, s'_t, t)$   
**Updates:**    `ensemble_buffer`        replay buffer of  $K$ -parallel perturbed data

```

1: for  $k$  in  $(1, \dots, K)$  do
2:   if  $m_t^k \sim \text{Unif}(\{0, 1\}) = 1$  then
3:     ensemble_buffer[k].enqueue((s_t, a_t, r_t, s'_t, t))

```

---

Algorithm 2 describes the `learn_from_buffer` method for the agent. For our experiments, we sometimes amend Algorithm 3 to learn periodically every  $N$  steps, rather than only at the end of the episode, but we mention this in the text where this is the case. This practice is common for most implementations of DQN and other reinforcement learning algorithms, but it does not play a significant role in our algorithm.

The final piece to describe BootDQN+prior is the `act` method for action selection. We employ a form of approximate Thompson sampling for RL via randomized value functions. Every episode, the agent selects  $j \sim \text{Unif}(1, \dots, K)$  and follows the greedy policy for  $Q_j$  for the duration of the episode.

## B Reinforcement learning experiments

In this section, we expand on details for the experimental set-up together with some additional results. Unless otherwise stated we use TensorFlow defaults, Adam optimizer with learning rate  $10^{-3}$  and uniform experience replay with batch size 128. For our  $\epsilon$ -greedy DQN baseline, we anneal epsilon linearly over 2000 episodes and perform hyperparameter sweeps over the initial epsilon  $\epsilon_0$ . All other agents (NoisyNet, Dropout, Ensemble, Bootstrap) use greedy policies according to an appropriate per-episode Thompson sampling.

### B.1 Chain environments

Figure 3 shows the time it takes each agent to learn a problem of size  $N$ . Figure 8 reproduces these results but on a log-log scale, which helps to reveal the problem scaling as  $N$  increases. As in Figure 3, the dashed line corresponds to a dithering lower bound  $T_{\text{learn}} = 2^N$ . We also include a solid line with slope equal to three, corresponding to a polynomial growth  $T_{\text{learn}} = \tilde{O}(N^3)$ .

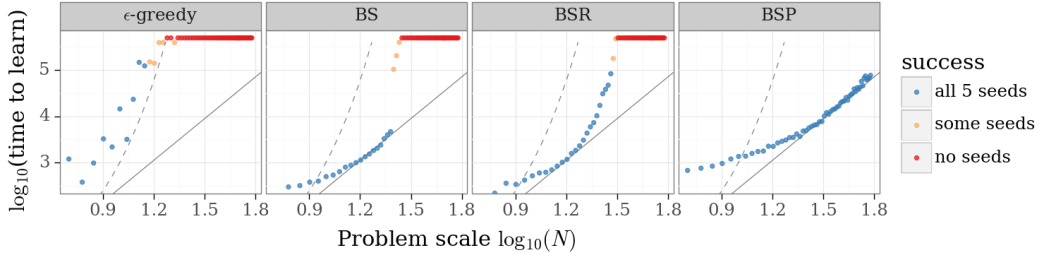


Figure 8: Log-log plot demonstrates scaling of learning behaviour.

In addition to BSP, BSR, BS and  $\epsilon$ -greedy displayed in Figure 3, we also ran parameter sweeps for dropout, NoisyNet and a count-based exploration strategy. Figure 9 presents the result for NoisyNet and dropout, each individually tuned up to 50k episodes. Even after tuning dropout rate and sampling frequency (by episode or by timestep) neither dropout nor NoisyNet scale successfully to large domains.

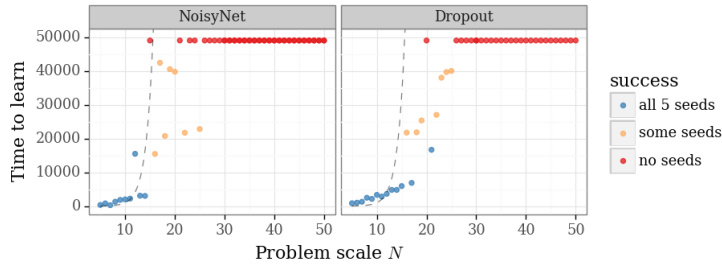


Figure 9: Learning time for noisy and dropout; neither approach scales well.

To compare with ‘count-based’ exploration we implement a version of DQN that optimizes the true reward plus a UCB exploration bonus  $\frac{\beta}{\sqrt{N_t(s)}}$ , where  $N(s)$  is the number of visits to state  $s$  prior to time  $t$  [26, 6]. Figure 10 shows that this count-based exploration strategy performs much worse than BSP, even after sweeping over bonus scale  $\beta$  and even with access to the true state visit-counts. This mirrors the outperformance of PSRL vs UCRL in tabular reinforcement learning. One explanation for this discrepancy comes from the inefficient way UCB-style algorithms propagate uncertainty over many timesteps [51, 44].

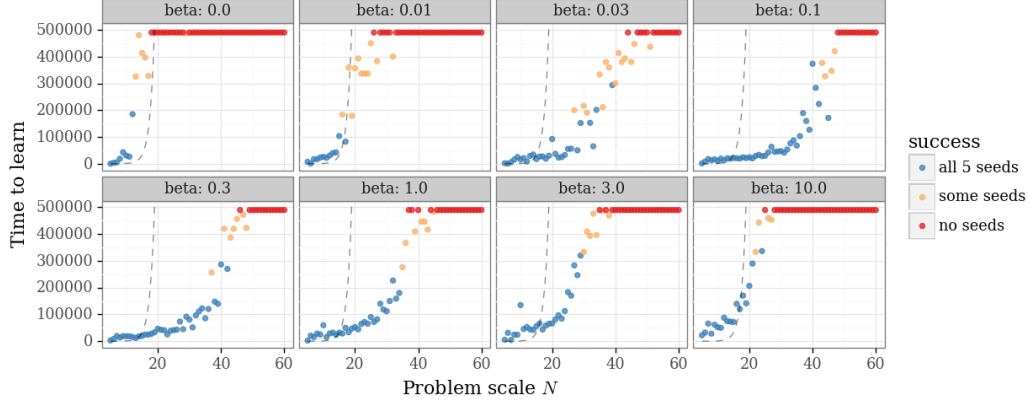


Figure 10: Sweeping over optimistic bonus; no scale of  $\beta$  matches BSP performance.

For all of our algorithms we tune agent hyperparameters by grid search. These were:

- **$\epsilon$ -greedy**:  $\epsilon = 0.1$ , linearly annealed to zero.
- **BSP**: prior scale  $\beta = 10$  (Figure 4b).
- **BSR**:  $l_2$  regularizer scale  $\lambda = 0.1$  (Figure 4a).
- **Dropout**: Resample mask every step with  $p_{\text{keep}} = 0.1$ .
- **NoisyNet**: Resample noise every step.
- **UCB**: Optimistic bonus  $\beta = 0.1$  (Figure 10).

## B.2 Sparse cartpole swing-up

In Section 4.2.2 we presented experiments showing that BSP outperforms benchmark algorithms. Figure 11 presents the sensitivity of BSP sensitivity to the prior scale  $\beta$  on this domain. Small values of  $\beta$  prematurely and suboptimally converge to the stationary policy, and so receive zero cumulative reward. Larger values of  $\beta$  take longer to wash away their prior effect, but we expect them to learn a performant policy eventually. This behaviour mirrors the scaling we saw in the chain environments, which is reassuring.

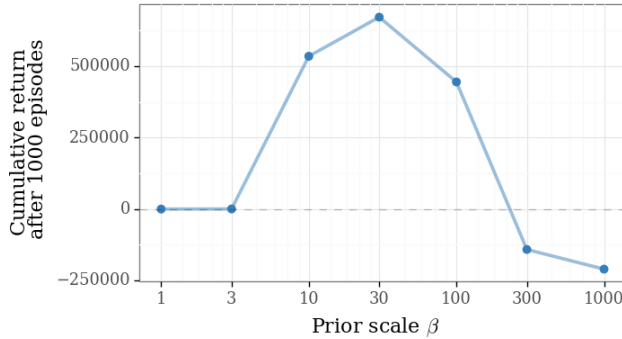


Figure 11: Sensitivity of performance to prior scale  $\beta$ .

## B.3 Montezuma’s revenge

In our experiments, we use the standard Atari configuration and preprocessing including greyscaling, frame stacking, action repeats, and random no-op starts [42], and the same agent hyperparameters as those used in the Ape-X paper [25]. However, our agent implementation is somewhat different and so our baseline results are not directly comparable across all games.



## C Why do we need a ‘prior’ mechanism for deep RL?

Section 2 outlines the need for a prior mechanism in deep RL, together with key failure cases for several of the most popular approaches. Due to space limitations we provide only simple illustrations of potential inadequacies of each method and this does not preclude their efficacy on any particular domain. In this appendix we expand on the details provided in Section 2 and provide suggestions for how these approaches might be remedied by future work.

### C.1 Dropout as posterior approximation

Previous work has suggested that dropout works as an effective variational approximation to the Bayesian posterior in neural networks, without special consideration for the network architecture [20]. However, Lemma 1 is a general statement that gives us cause to question the quality of this approximation. In this subsection we dig deeper into an extremely simple estimation problem, a linear network with  $d$  units to estimate the mean of a random variable  $Y \in \mathbb{R}$ . Even in this simple setting dropout performs poorly as a Bayesian posterior.

We form predictions  $f_\theta = \sum_{i=1}^d w_i \theta_i$  with  $w_i \sim \text{Ber}(p)$ , square loss and regularizer  $\mathcal{R}(\theta) = \lambda \|\theta\|^2$  for  $\lambda > 0$ . Then for any data  $\mathcal{D}$  with empirical mean  $\bar{y}$ , the expected loss solution to (1) is given by [61]  $\theta_p^* = \bar{\theta} \mathbf{1}$  for  $\bar{\theta} = \frac{\bar{y}}{1+p(d-1)+\frac{\lambda}{d}}$ .<sup>6</sup> The resultant predictive distribution therefore has mean  $\mu = \bar{\theta} dp$  and standard deviation  $\sigma = \bar{\theta} \sqrt{dp(1-p)}$ .

If we are to understand dropout as an approximation to a Bayesian posterior, then we should note that this behavior is unusual. First, the only connection to the data is through the empirical mean  $\bar{y}$ ; any possible dataset with the same mean would result in the same ‘posterior’ distribution. Second we note that  $\sigma = \mu \sqrt{(1-p)/dp}$ . This coupling means it is not possible for  $\sigma \rightarrow 0$  and  $\mu \rightarrow 0$ , regardless of  $\lambda$ . More typically we would imagine  $\mu \rightarrow \mathbb{E}[Y]$  and  $\sigma \rightarrow 0$  according to the Bayesian central limit theorem [12].

This disconnect is not simply an analytical mistake, but can lead to arbitrarily bad decisions in even the simplest problems. Imagine a simple two-armed bandit problem with one arm’s rewards  $\sim \text{Ber}(1/2)$  and the other’s  $\sim \text{Ber}(1/2+\epsilon)$ , and the agent does not know which arm is which a priori. This style of problem is particularly well understood with guarantees that Thompson sampling with more reasonable forms of posterior approximation incur regret  $\tilde{O}(\log(T))$  in this setting [2]. We refer to this problem as  $\dagger$ . The following result highlights that dropout as posterior approximation can perform poorly even on this simple domain.

**Lemma 4** (Dropout sampling attains linear regret in  $\dagger$ ).

*Fix any  $d, p, \lambda$  and consider the problem of  $\dagger$  with an agent employing Thompson sampling by dropout for action selection. Then the expected regret is  $\Omega(T)$ .*

*Proof.* For any  $d, p, \lambda$  and any observed data  $\mathcal{H}_t$ , there exists a non-zero probability  $P_1(s, p, \lambda, \mathcal{H}_t) > \frac{p^d}{2}$  of selecting action 1 over action action 2. We can see this by imagining all units estimating action 2 are set to zero, then there is at least 50% chance of selection action 1. This proves<sup>7</sup> that  $\mathbb{E}[\text{Regret}(T)] \geq \frac{ep^d}{2} T$  for all  $T$ .  $\square$

Although our analysis of dropout has focused on an exceedingly simple functional form, the key insight that the degree of variability in the posterior distribution does not concentrate with data extends to any neural network architecture. Figure 12 presents the dropout posterior on a simple regression task with a (20,20)-MLP with rectified linear units. We display the predictive distribution under varying amounts of data. The dropout sampling distribution does not converge with increasing amounts of data, whereas the bootstrapped sampling approach behaves much more reasonably. This leads to poor performance in reinforcement learning tasks too, as we saw in Appendix B.

<sup>6</sup>This corrects an errant derivation in [46], but maintains the same overall message.

<sup>7</sup>Note that this lower bound is very conservative and provided only for illustration. A more precise analysis would show poor performance even for large  $d$ .

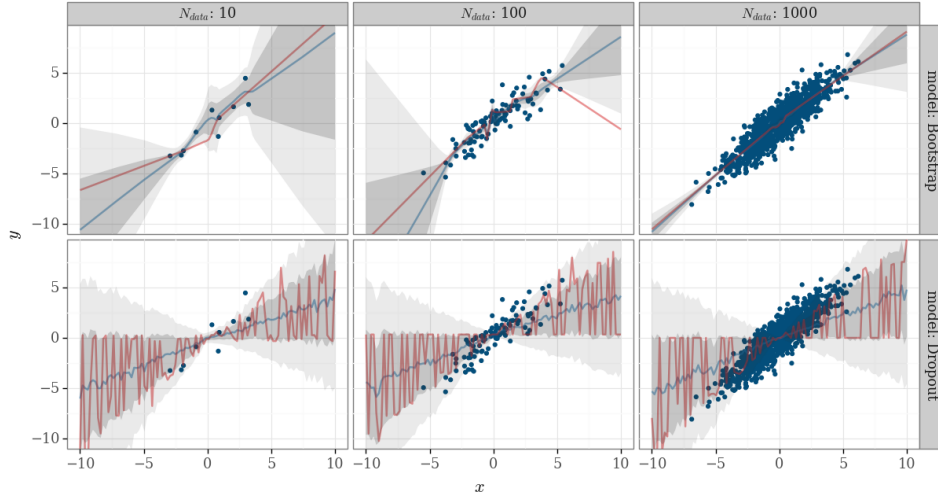


Figure 12: Dropout does not converge with increasing data even with a complex neural network. Grey regions indicate  $\pm 1, 2$  standard deviations, the mean is shown in blue and a single posterior sample in red.

## C.2 Variational inference and Bellman error

Lemma 2 highlights that the basic loss most commonly used by variational approximations to the value function are fundamentally ill-suited to the problem at hand [35, 17]. In Appendix B we present results of such a variational approach, NoisyNet, to some of our benchmark reinforcement learning tasks. As expected, the algorithm performs poorly even after extensive tuning. At the heart of this issue is a sample-based loss that trains to match the *expectation* of the target distribution, but does not attempt to match the higher moments of the uncertainty. However, we could imagine an alternative approach that does aim to match the entire resultant distribution, for example via parameterized distribution and cross entropy loss; we leave this to future work.

Even where variational inference is employed correctly over multiple timesteps, it may be difficult to encode useful prior knowledge in common variational methods. First, many applications of variational inference (VI) model the distribution over network weights as a product of independent Gaussians [11]. These models facilitate efficient computation, but can underestimate the uncertainty and may be a poor choice for encoding prior knowledge. Even if one were given a mapping of prior knowledge to weights, the confounding demands of good initialization for SGD training may interfere negatively [23]. For this reason practical applications of VI to RL typically use very little prior effect, or even no prior regularization at all [35, 17].

## C.3 ‘Distributional reinforcement learning’

Unlike the objections of Appendix C.2, ‘distributional RL’<sup>8</sup> does learn a value function estimate through a distributional loss. However, this distribution is a distribution over *outcomes* and not a distribution over the *epistemic uncertainty* in the mean beliefs. This distinction between two types of uncertainty, (1.) things that you don’t know and (2.) things that are stochastic, is a delicate one and is important to characterize correctly. Both are discussed under many names:

- (1.) ‘Reducible uncertainty’  $\iff$  ‘epistemic uncertainty’  $\iff$  ‘uncertainty’,
- (2.) ‘Irreducible uncertainty’  $\iff$  ‘aleatoric uncertainty’  $\iff$  ‘risk’.

<sup>8</sup>Any method for Bayesian RL might reasonably claim to be a distributional perspective on reinforcement learning. For this reason, we use quotation marks when we want to distinguish the specific form of distributional RL popularized by [6].

Typical decision problems may include elements of both types of uncertainty. Flipping a coin we might want to know both (1.) our posterior beliefs over the probability of heads and (2.) a distribution that categorizes the likely possible outcomes. However, it should be clear that the two concepts are fundamentally distinct. For the purposes of exploration, the Bayesian uncertainty over (1.) should prioritize the acquisition of new knowledge. ‘Distributional RL’ approximates (2.) and its role is not exchangeable with (1.).

**Lemma 5** (Using ‘distributional RL’ as a posterior can lead to arbitrarily bad decisions). *Consider an agent with full information that decides between action 1 with reward  $\sim \text{Ber}(0.5)$  and action 2 with reward  $1-\epsilon$  for  $0 < \epsilon < \frac{1}{2}$ . If the agent employs Thompson sampling correctly then it will pick action 2 at every step with zero regret. If the agent mistakenly employs Thompson sampling over its ‘distributional value function’ then it will incur*

$$\mathbb{E}[\text{Regret}(T)] \geq \frac{1}{2} \left( \frac{1}{2} - \epsilon \right) T.$$

Lemma 5 shows that using the ‘distributional’ value function approximating (2.) can be a poor proxy for the Bayesian uncertainty. However, the Bayesian uncertainty can be a similarly poor proxy for the ‘distributional’ value function. This can be equally damaging, particularly if the agent has some risk-sensitive utility with respect to cumulative rewards. It is entirely possible to combine both notions of uncertainty in an agent, although for the goal of maximizing expected cumulative it is not entirely clear what is the benefit of modeling (2.). Certainly, ‘distributional’ agents have recently attained strong scores in Atari 2600 benchmarks but it is so far unclear exactly what the source of this outperformance comes from [8, 13]. Possible explanations may include more stable gradients, bounded values and the ‘many predictions’ hypothesis [64]: that learning a distribution may effectively create a series of auxiliary losses. We leave these questions for future work.

#### C.4 Count-based uncertainty estimates

Count-based approaches to exploration give a bonus to states that have not been visited frequently according to some density measure  $p(x)$ . These methods have performed well in many sparse reward tasks such as Montezuma’s Revenge, where visiting new states acts as a shaping reward for the true reward [6]. However, a count-based bonus is generally a poor approach to exploration beyond the tabular setting. To see why this is the case note that the density measure of the states may not correlate well with the agent’s uncertainty over the optimal policy in that state. We can imagine situations both where the state is visually new, but an agent should still know exactly what to do; and also settings that are only delicately different to a common situation but still necessitate exploration of the optimal policy.

This disconnect shows up in problems as simple as the linear bandit.<sup>9</sup> Via a packing argument, an agent with count-based uncertainty will require  $\tilde{O}\left(\frac{1}{\epsilon^d}\right)$  measurements to cover the space up to radius  $\epsilon$ . By contrast, an agent that explores this space efficiently can resolve its uncertainty in only  $\tilde{O}(d)$  measurements. Thompson sampling with a linear model naturally recovers this performance [57]. The fact that this failure can arise even in a linear system, and even when the density can be estimated precisely, suggests that count-based exploration is not *in general* an effective method for simultaneous exploration with generalization; even if it may be effective at some specific tasks.

#### C.5 Ensembles without priors

This paper builds upon a line of research that uses an ensemble of trained models to approximate a posterior distribution. Compared to previous works, our main contribution is to highlight the importance of a ‘prior’ mechanism in ensemble uncertainty. Figure 13 presents an extremely simple example of 1D regression with a (20,20)-MLP and rectified linear units. The data consists of  $x_i = \frac{i-5}{5}$  for  $i=0, \dots, 10$  and  $y_i = 51\{i=10\}$ .

The results above highlight the drawbacks of naive ensembles. A pure ensemble trained from random initializations fits the data exactly and leads to almost zero uncertainty anywhere in

<sup>9</sup>Reward  $r_t(x_t) = x_t^T \theta^* + \epsilon_t$  for some  $\theta^* \in \mathbb{R}^d$  and  $\epsilon_t \sim N(0,1)$  [56].

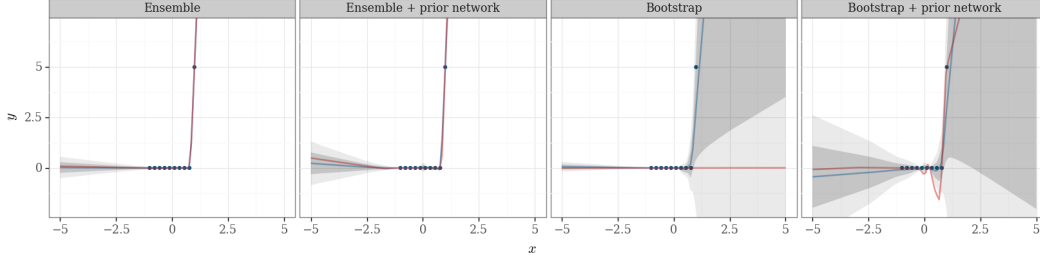


Figure 13: Posterior predictive distributions for ensemble uncertainty.

the space [31]. A bootstrapped ensemble takes the variability of the data into account and thus has a wide predictive uncertainty as  $x$  grows large and positive. However, where the data has target value zero, bootstrapping will always produce a zero target and consequently the ensemble has almost zero predictive uncertainty as  $x$  becomes large and negative [47].

This lack of prior uncertainty can lead to arbitrarily poor decisions, as outlined in [50]. If an agent has only ever observed zero reward, then no amount of bootstrapping or ensembling will cause it to simulate positive rewards. This issue is easily remedied by the addition of a prior mechanism, either through  $l_2$  regularization to initial random weights (4), or the addition of a fixed additive random ‘prior network’ (5).

## C.6 Summary

We summarize the issues raised in Section 2 in Table 1. This table is meant only as a rough summary and should not be taken as rigorous statement. Roughly speaking, a green tick means success, red cross means failure and a yellow circle means something in between. This paper proposes a combination of bootstrap sampling with prior function as an effective computational approximation to Bayesian inference in deep RL. Although our method is somewhat computationally expensive, since it requires training an ensemble of models instead of one, this computation can be done in parallel and so is amenable to large scale distributed computation.

Table 1: Important issues in posterior approximations for deep reinforcement learning.

	Data conc.	Learned metric	Multi step	Works in noise	Prior effect	Cheap compute
Dropout [20]	✗	✓	✗	●	✗	✓
NoisyNet [17]	●	✓	✗	✓	✗	✓
BBB / VI [11]	●	✓	✗	✓	●	✓
Density count [6]	✓	✗	✓	✗	●	✓
‘Distributional’ RL [8]	✗	●	✓	●	✗	✓
Ensemble [31]	✗	✓	✓	✗	✗	✗
Bootstrap [47]	✓	✓	✓	✓	✗	✗
<b>Bootstrap + prior</b>	✓	✓	✓	✓	✓	✗
Exact Bayes	✓	✓	✓	✓	✓	✗✗✗