

State-Aware Variational Thompson Sampling for Deep Q-Networks

Siddharth Aravindan

National University of Singapore
siddharth.aravindan@comp.nus.edu.sg

Wee Sun Lee

National University of Singapore
leews@comp.nus.edu.sg

ABSTRACT

Thompson sampling is a well-known approach for balancing exploration and exploitation in reinforcement learning. It requires the posterior distribution of value-action functions to be maintained; this is generally intractable for tasks that have a high dimensional state-action space. We derive a variational Thompson sampling approximation for DQNs which uses a deep network whose parameters are perturbed by a learned variational noise distribution. We interpret the successful NoisyNets method [10] as an approximation to the variational Thompson sampling method that we derive. Further, we propose State Aware Noisy Exploration (SANE) which seeks to improve on NoisyNets by allowing a non-uniform perturbation, where the amount of parameter perturbation is conditioned on the state of the agent. This is done with the help of an auxiliary perturbation module, whose output is state dependent and is learnt end to end with gradient descent. We hypothesize that such state-aware noisy exploration is particularly useful in problems where exploration in certain *high risk* states may result in the agent failing badly. We demonstrate the effectiveness of the state-aware exploration method in the off-policy setting by augmenting DQNs with the auxiliary perturbation module.

KEYWORDS

Deep Reinforcement Learning; Thompson Sampling; Exploration

ACM Reference Format:

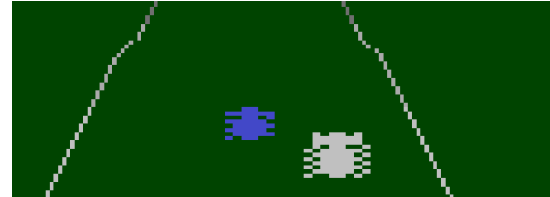
Siddharth Aravindan and Wee Sun Lee. 2021. State-Aware Variational Thompson Sampling for Deep Q-Networks. In *Proc. of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021)*, Online, May 3–7, 2021, IFAAMAS, 12 pages.

1 INTRODUCTION

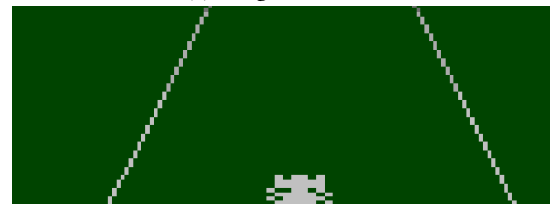
Exploration is a vital ingredient in reinforcement learning algorithms that has largely contributed to its success in various applications [12, 15–17]. Traditionally, deep reinforcement learning algorithms have used naive exploration strategies such as ϵ -greedy, Boltzmann exploration or action-space noise injection to drive the agent towards unfamiliar situations. Although effective in simple tasks, such undirected exploration strategies do not perform well in tasks with high dimensional state-action spaces.

Theoretically, Bayesian approaches like Thompson sampling have been known to achieve an optimal exploration-exploitation trade-off in multi-armed bandits [2, 3, 14] and also have been shown to provide near optimal regret bounds for Markov Decision Processes (MDPs) [4, 21]. Practical usage of such methods, however, is

Proc. of the 20th International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2021), U. Endriss, A. Nowé, F. Dignum, A. Lomuscio (eds.), May 3–7, 2021, Online. © 2021 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org). All rights reserved.



(a) A high risk state



(b) A low risk state

Figure 1: The white car which is controlled by the agent, has to move forward while avoiding other cars. (a) In this state, any action other than moving straight will result in a crash, making it a high risk state. (b) This is a low risk state since exploring random actions will not lead to a crash.

generally intractable as they require the posterior distribution of value-action functions to be maintained.

At the same time, in practical applications, perturbing the parameters of the model with Gaussian noise to induce exploratory behaviour has been shown to be more effective than ϵ -greedy and other approaches that explore primarily by randomization of the action space [10, 24]. Furthermore, adding noise to model parameters is relatively easy and introduces minimal computational overhead. NoisyNets [10], in particular, has been known to achieve better scores on the full Atari suite than other exploration techniques[31].

In this paper, we derive a variational Thompson sampling approximation for Deep Q-Networks (DQNs), where the model parameters are perturbed by a learned variational noise distribution. This enables us to interpret NoisyNets as an approximation of Thompson sampling, where minimizing the NoisyNet objective is equivalent to optimizing the variational objective with Gaussian prior and approximate posterior distributions. These Gaussian approximating distributions, however, apply perturbations uniformly across the agent’s state space. We seek to improve this by approximating the Thompson sampling posterior with Gaussian distributions whose variance is dependent on the agent’s state.

To this end, we propose State Aware Noisy Exploration (SANE), an exploration strategy that induces exploration through state dependent parameter space perturbations. These perturbations are

added with the help of an augmented state aware perturbation module, which is trained end-to-end along with the parameters of the main network by gradient descent.

We hypothesize that adding such perturbations helps us mitigate the effects of high risk state exploration, while exploring effectively in low risk states. We define a high risk state as a state where a wrong action might result in adverse implications, resulting in an immediate failure or transition to states from which the agent is eventually bound to fail. Exploration in such states might result in trajectories similar to the ones experienced by the agent as a result of past failures, thus resulting in low information gain. Moreover, it may also prevent meaningful exploratory behaviour at subsequent states in the episode, that may have been possible had the agent taken the correct action at the state. A low risk state, on the other hand, is defined as a state where a random exploratory action does not have a significant impact on the outcome or the total reward accumulated by the agent within the same episode. A uniform perturbation scheme for the entire state space may thus be undesirable in cases where the agent might encounter high risk states and low risk states within the same task. An instance of a high risk state and low risk state in Enduro, an Atari game, is shown in Figure 1. We try to induce uncertainty in actions, only in states where such uncertainty is needed through the addition of a state aware perturbation module.

To test our assumptions, we experimentally compare two SANE augmented Deep Q-Network (DQN) variants, namely the simple-SANE DQN and the Q-SANE DQN, with their NoisyNet counterparts [10] on a suite of 11 Atari games. Eight of the games in the suite have been selected to have high risk and low risk states as illustrated in Figure 1, while the remaining three games do not exhibit such properties. We find that agents that incorporate SANE do better in most of the eight games. An added advantage of SANE over NoisyNets is that it is more scalable to larger network models. The exploration mechanism in NoisyNets [10] adds an extra learnable parameter for every weight to be perturbed by noise injection, thus tying the number of parameters in the exploration mechanism to the network architecture being perturbed. The noise-injection mechanism in SANE on the other hand, is a separate network module, independent of the architecture being perturbed. The architecture of this perturbation module can be modified to suit the task. This makes it more scalable to larger networks.

2 BACKGROUND

2.1 Markov Decision Processes

A popular approach towards solving sequential decision making tasks involves modelling them as MDPs. A MDP can be described as a 5-tuple, $(\mathcal{S}, \mathcal{A}, \mathcal{R}(\cdot), \mathcal{T}(\cdot, \cdot), \gamma)$, where \mathcal{S} and \mathcal{A} denote the state space and action space of the task, \mathcal{T} and \mathcal{R} represent the state-transition and reward functions of the environment respectively and γ is the discount factor of the MDP. Solving a MDP entails learning an optimal policy that maximizes the expected cumulative discounted reward accrued during the course of an episode. Planning algorithms can be used to solve for optimal policies, when \mathcal{T} and \mathcal{R} are known. However, when these functions are unavailable, reinforcement learning methods help the agent *learn* good policies.

2.2 Deep Q-Networks

A DQN [17] is a value based temporal difference learning algorithm, that estimates the action-value function by minimizing the temporal difference between two successive predictions. It uses a deep neural network as a function approximator to compute all action values of the optimal policy $Q^{\pi^*}(s, a)$, for a given a state s . A typical DQN comprises two separate networks, the Q network and the target network. The Q network aids the agent in interacting with the environment and collecting training samples to be added to the experience replay buffer, while the target network helps in calculating *target* estimates of the action value function. The network parameters are learned by minimizing the loss $\mathcal{L}(\theta)$ given in Equation 1, where θ and θ' are the parameters of the Q network and the target network respectively. The training instances (s_i, a_i, r_i, s'_i) are sampled uniformly from the experience replay buffer, which contains the k most recent transitions experienced by the agent.

$$\mathcal{L}(\theta) = \mathbb{E} \left[\frac{1}{b} \sum_{i=1}^b (Q(s_i, a_i; \theta) - (r_i + \gamma \max_a Q(s'_i, a; \theta')))^2 \right] \quad (1)$$

2.3 Thompson Sampling

Thompson sampling [32] works under the Bayesian framework to provide a well balanced exploration-exploitation trade-off. It begins with a prior distribution over the action-value and/or the environment and reward models. The posterior distribution over these models/values is updated based on the agent’s interactions with the environment. A Thompson sampling agent tries to maximize its expected value by acting greedily with respect to a sample drawn from the posterior distribution. Thompson sampling has been known to achieve optimal and near optimal regret bounds in stochastic bandits [2, 3, 14] and MDPs [4, 21] respectively.

3 RELATED WORK

Popularly used exploration strategies like ϵ -greedy exploration, Boltzmann exploration and entropy regularization [30], though effective, can be wasteful at times, as they do not consider the agent’s uncertainty estimates about the state. In tabular settings, count based reinforcement learning algorithms such as UCRL [5, 13] handle this by maintaining state-action visit counts and incentivize agents with exploration bonuses to take actions that the agent is uncertain about. An alternative approach is suggested by posterior sampling algorithms like PSRL [29], which maintain a posterior distribution over the possible environment models, and act optimally with respect to the model sampled from it. Both count based and posterior sampling algorithms have convergence guarantees in this setting and have been proven to achieve near optimal exploration-exploitation trade-off. Unfortunately, sampling from a posterior over environment models or maintaining visit counts in most real world applications are computationally infeasible due to the high dimensional state space and action space involved with these tasks. However, approximations of such methods that do well have been proposed in recent times.

Bellemare et al. [7], generalizes count based techniques to non-tabular settings by using pseudo-counts obtained from a density model of the state space, while [28] follows a similar approach but

uses a predictive model to derive the bonuses. Ostrovski et al. [23] builds upon [7] by improving upon the density models used for granting exploration bonuses. Additionally, surprise-based motivation [1] learns the transition dynamics of the task, and adds a reward bonus proportional to the Kullback–Leibler (KL) divergence between the true transition probabilities and the learned model to capture the agent’s *surprise* on experiencing a transition not conforming to its learned model. Such methods that add exploration bonuses prove to be most effective in settings where the rewards are very sparse but are often complex to implement [24].

Randomized least-squares value iteration (RLSVI) [22] is an approximation of posterior sampling approaches to the function approximation regime. RLSVI draws samples from a distribution of linearly parameterized value functions, and acts according to the function sampled. [19] and [20] are similar in principle to [22]; however, instead of explicitly maintaining a posterior distribution, samples are procured with the help of bootstrap re-sampling. Randomized Prior Functions [18] adds *untrainable prior networks* with the aim of capturing uncertainties not available from the agent’s experience, while [6] tries to do away with duplicate networks by using Bayesian linear regression with Gaussian prior. Even though the action-value functions in these methods are no longer restricted to be linear, maintaining a bootstrap or computing a Bayesian linear regression makes these methods computationally expensive compared to others.

Parameter perturbations which form another class of exploration techniques, have been known to enhance the exploration capabilities of agents in complex tasks [9, 24, 33]. Rückstieß et al. [25] show that this type of policy perturbation in the parameter space outperforms action perturbation in policy gradient methods, where the policy is approximated with a linear function. However, Rückstieß et al. [25] evaluate this on tasks with low dimensional state spaces. When extended to high dimensional state spaces, black box parameter perturbations [26], although proven effective, take a long time to learn good policies due to their non adaptive nature and inability to use gradient information. Gradient based methods that rely on adaptive scaling of the perturbations, drawn from spherical Gaussian distributions [24], gradient based methods that learn the amount of noise to be added [10] and gradient based methods that learn dropout policies for exploration [33] are known to be more sample efficient than black box techniques. NoisyNets [10], a method in this class, has been known to demonstrate consistent improvements over ϵ -greedy across the Atari game suite unlike other count-based methods [31]. Moreover, these methods are also often easier to implement and computationally less taxing than the other two classes of algorithms mentioned above.

Our exploration strategy belongs to the class of methods that perturb parameters to effect exploration. Our method has commonalities with the parameter perturbing methods above as we sample perturbations from a spherical Gaussian distribution whose variance is learnt as a parameter of the network. However, the variance learnt, unlike NoisyNets [10], is conditioned on the current state of the agent. This enables it to sample perturbations from different Gaussian distributions to vary the amount of exploration when the states differ. Our networks also differ in the type of perturbations applied to the parameters. While [10] obtains a noise sample from possibly different Gaussian distributions for each parameter, our

network, like [24], samples all perturbations from the same, but state aware, Gaussian distribution. Moreover, the noise injection mechanism in SANE is a separate network module that is subject to user design. This added flexibility might make it more scalable to larger network models when compared to NoisyNets, where this mechanism is tied to the network being perturbed.

4 VARIATIONAL THOMPSON SAMPLING

Bayesian methods like Thompson Sampling use a posterior distribution $p(\theta|\mathcal{D})$ to sample the weights of the neural network, given \mathcal{D} , the experience collected by the agent. $p(\theta|\mathcal{D})$ is generally intractable to compute and is usually approximated with a variational distribution $q(\theta)$. Let $D = (X, Y)$ be the dataset on which the agent is trained, with X being the set of inputs, and Y being the target labels. Variational methods minimize the KL divergence between $q(\theta)$ and $p(\theta|D)$ to make $q(\theta)$ a better approximation. Appendix A shows that minimizing $KL(q(\theta), p(\theta|D))$ is equivalent to maximizing the Evidence Lower Bound (ELBO), given by Equation 2.

$$ELBO = \int q(\theta) \log p(Y|X, \theta) d\theta - KL(q(\theta), p(\theta)) \quad (2)$$

For a dataset with N datapoints, and under the i.i.d assumption, we have :

$$\log p(Y|X, \theta) = \log \prod_{i=1}^N p(y_i|x_i, \theta) = \sum_{i=1}^N \log p(y_i|x_i, \theta)$$

So, the objective to maximize is :

$$\begin{aligned} & \max \int q(\theta) \sum_{i=1}^N \log p(y_i|x_i, \theta) d\theta - KL(q(\theta), p(\theta)) \\ & = \max \left[\left(\sum_{i=1}^N \int q(\theta) \log p(y_i|x_i, \theta) d\theta \right) - KL(q(\theta), p(\theta)) \right] \end{aligned}$$

In DQNs, the inputs x_i are state-action tuples, and the its corresponding target y_i is an estimate of $Q(s, a)$. Traditionally, DQNs are trained by minimizing the squared error, which assumes a Gaussian error distribution around the target value. Assuming the same, we define $\log p(y_i|x_i, \theta)$ in Equation 3, where y_i is the approximate target Q value of (s_t^i, a_t^i) given by $T_i = r_t^i + \max_a \gamma Q(s_{t+1}^i, a; \theta)$, σ_ϵ^2 is the variance of the error distribution and $C(\sigma_\epsilon) = -\log \sqrt{(2\pi)\sigma_\epsilon}$.

$$\log p(y_i|x_i, \theta) = \frac{-(Q(s_t^i, a_t^i; \theta) - T_i)^2}{2\sigma_\epsilon^2} + C(\sigma_\epsilon) \quad (3)$$

$$\begin{aligned} ELBO &= \left(\sum_{i=1}^N \int q(\theta) \left[\frac{-(Q(s_t^i, a_t^i; \theta) - T_i)^2}{2\sigma_\epsilon^2} \right] d\theta \right) \\ &\quad - KL(q(\theta), p(\theta)) + NC(\sigma_\epsilon) \\ &= C_1 \left(\sum_{i=1}^N \int q(\theta) [-(Q(s_t^i, a_t^i; \theta) - T_i)^2] d\theta \right) \\ &\quad - KL(q(\theta), p(\theta)) + NC(\sigma_\epsilon) \end{aligned} \quad (4)$$

We approximate the integral for each example with a Monte Carlo estimate by sampling a $\hat{\theta}_i \sim q(\theta)$, giving

$$ELBO \approx C_1 \left(\sum_{i=1}^N -(Q(s_t^i, a_t^i; \hat{\theta}_i) - T_i)^2 \right) - KL(q(\theta), p(\theta)) + NC(\sigma_\epsilon)$$

As $C(\sigma_\epsilon)$ is a constant with respect to θ , maximizing the ELBO is approximately the same as optimizing the following objective.

$$\max \left[C_1 \left(\sum_{i=1}^N -(Q(s_t^i, a_t^i; \hat{\theta}_i) - T_i)^2 \right) - KL(q(\theta), p(\theta)) \right] \quad (5)$$

4.1 Variational View of NoisyNet DQNs

The network architecture of NoisyNet DQNs usually comprises a series of convolutional layers followed by some fully connected layers. The parameters of the convolutional layers are not perturbed, while every parameter of the fully connected layers is perturbed by a separate Gaussian noise whose variance is learned along with the other parameters of the network.

For the unperturbed parameters of the convolutional layers, we consider $q(\theta_c) = \mathcal{N}(\mu_c, \epsilon I)$. The parameters of any neural network are usually used in the floating point format. We choose a value of ϵ that is close enough to zero, such that adding any noise sampled from these distributions does not change the value of the weight as represented in this format with high probability. For the parameters of the fully connected layers, we take $q(\theta_{fc}) = \mathcal{N}(\mu_{fc}, \Sigma_{fc})$ where Σ is a diagonal matrix with Σ_{fc}^{ii} equal to the learned variance for the parameter θ_i . We take the prior $p(\theta) = \mathcal{N}(0, I)$ for all the parameters of the network.

With this choice of $p(\theta)$ and $q(\theta)$, the value of $KL(q(\theta), p(\theta))$ can be computed as shown in Equation 6, where k_1 and k_2 are the number of parameters in the convolutional and fully connected layers respectively. Note that k_1 , k_2 and ϵ are constants given the network architecture.

$$KL(q(\theta), p(\theta)) = \frac{1}{2} [-k_1 \log(\epsilon) + k_1 \epsilon + \|\mu_c\|_2^2 - k_1] + \frac{1}{2} [-\log |\Sigma_{fc}| + \text{tr}(\Sigma_{fc}) + \|\mu_{fc}\|_2^2 - k_2] \quad (6)$$

As NoisyNet DQN agents are usually trained on several million interactions, we assume that the KL term is dominated by the log likelihood term in the ELBO. Thus, maximizing the objective in Equation (5) can be approximated by optimizing the following objective :

$$\max \left(\sum_{i=1}^N -(Q(s_t^i, a_t^i; \hat{\theta}_i) - T_i)^2 \right) \quad (7)$$

which is the objective that NoisyNet DQN agents optimize. In NoisyNets, every sample $\hat{\theta}_i \sim \mathcal{N}(\mu, \Sigma)$ is obtained by a simple reparameterization of the network parameters : $\hat{\theta}_i = \mu + \Sigma \epsilon$, where $\epsilon \sim \mathcal{N}(0, I)$. This reparameterization helps NoisyNet DQNs to learn through a sampled $\hat{\theta}_i$.

4.2 State Aware Approximating Distributions

It can be seen that the approximate posterior distribution $q(\theta)$ is state agnostic, i.e., it applies perturbations uniformly across the

state space, irrespective of whether the state is *high risk* or *low risk*. We thus postulate that $q(\theta|s)$ is potentially a better variational approximator. $q(\theta)$ is a special case of a state aware variational approximator where $q(\theta|s)$ is the same for all s . A reasonable ELBO estimate for such an approximate distribution would be to extend the ELBO in Equation 4 to accommodate $q(\theta|s)$ as shown in 8.

$$ELBO = C_1 \left(\sum_{i=1}^N \int q(\theta|s_t^i) [-(Q(s_t^i, a_t^i; \theta) - T_i)^2] d\theta \right) - \frac{1}{N} \sum_{i=1}^N KL(q(\theta|s_t^i), p(\theta)) + NC(\sigma_\epsilon) \quad (8)$$

Approximating the integral for each example with a Monte Carlo estimate by sampling a $\hat{\theta}_i \sim q(\theta|s_t^i)$, maximizing the ELBO is equivalent to solving 9.

$$\max \left[\sum_{i=1}^N \left(-C_1 (Q(s_t^i, a_t^i; \hat{\theta}_i) - T_i)^2 - \frac{1}{N} KL(q(\theta|s_t^i), p(\theta)) \right) \right] \quad (9)$$

We assume that the KL term will eventually be dominated by the log likelihood term in the ELBO, given a sufficiently large dataset. This posterior approximation leads us to the formulation of SANE DQNs as described in the following sections.

5 STATE AWARE NOISY EXPLORATION

State Aware Noisy Exploration (SANE), is a parameter perturbation based exploration strategy which induces exploratory behaviour in the agent by adding noise to the parameters of the network. The noise samples are drawn from the Gaussian distribution $\mathcal{N}(0, \sigma^2(h(s; \theta); \Theta))$, where $\sigma(h(s; \theta); \Theta)$ is computed as a function of a hidden representation, $h(s; \theta)$, of the state s of the agent by an auxiliary neural network module, i.e., $\sigma(h(s; \theta); \Theta) = g_\Theta(h(s; \theta))$, where Θ and θ refer to the parameters of the auxiliary perturbation network (g) and the parameters of the main network respectively.

5.1 State Aware Noise Sampling

To procure state aware noise samples, we first need to compute $\sigma(h(s; \theta); \Theta)$, the state dependent standard deviation of the Normal distribution from which the perturbations are sampled. As stated above, we do this by adding an auxiliary neural network module. $\sigma(h(s; \theta); \Theta)$ is then used to generate perturbations $\epsilon \sim \mathcal{N}(0, \sigma^2(h(s; \theta); \Theta))$ for every network parameter using noise samples from the standard Normal distribution, $\hat{\epsilon} \sim \mathcal{N}(0, 1)$, in tandem with a simple reparameterization of the sampling network [10, 24, 26] as shown in Equation 10.

$$\epsilon = \sigma(h(s; \theta); \Theta) \hat{\epsilon}; \quad \hat{\epsilon} \sim \mathcal{N}(0, 1) \quad (10)$$

State aware perturbations can be added to all types of layers in the network. The standard baseline architectures used by popular deep reinforcement learning algorithms for tasks such as playing Atari games mainly consist of several convolutional layers followed by multiple fully connected layers. We pass the output of the last convolutional layer as the hidden representation $h(s; \theta)$ to compute the state aware standard deviation, $\sigma(h(s; \theta); \Theta)$, where θ is the set of parameters of the convolutional layers. Perturbations using $\sigma(h(s; \theta); \Theta)$ are then applied to the following fully connected layers.

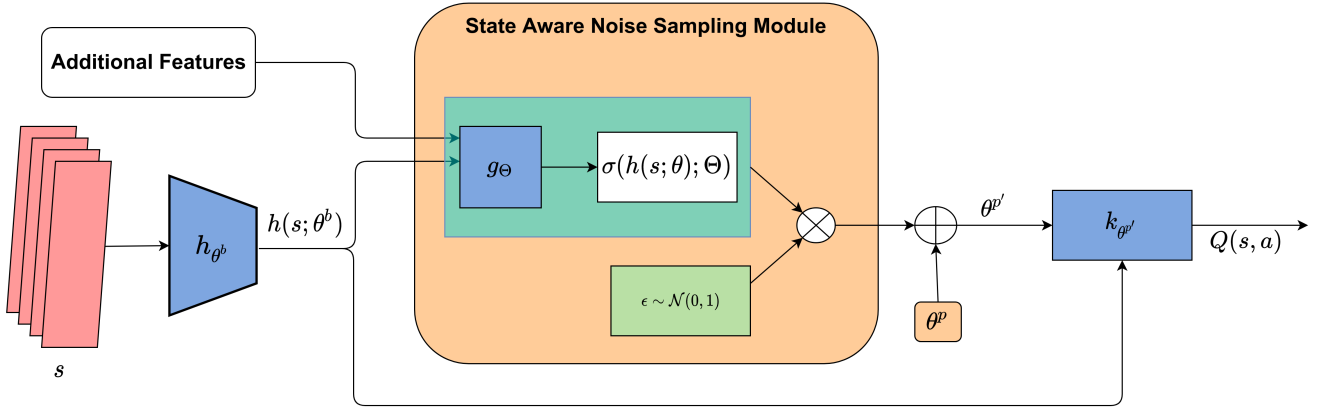


Figure 2: A high level view of a State Aware Noisy Exploring Network.

Our mechanism of introducing perturbations is similar to Noisy DQNs [10] and adaptive parameter space noise [24]. Given a vector $x \in \mathbb{R}^l$ as input to a fully connected layer with m outputs, an unperturbed layer computes a matrix transformation of the form $y = Wx + b$, where W and b are the parameters associated with the layer, and $y \in \mathbb{R}^m$. We modify such layers with state-aware perturbations, by adding noise elements sampled from $\mathcal{N}(0, \sigma^2(h(s; \theta); \Theta))$ (Equation 10). This results in the perturbed fully connected layer computing a transform equivalent to Equation 11, where $\tilde{W} = W + \sigma(h(s; \theta); \Theta)\epsilon_w$, $\tilde{b} = b + \sigma(h(s; \theta); \Theta)\epsilon_b$, and $\epsilon_w \in \mathbb{R}^{m \times l}$, $\epsilon_b \in \mathbb{R}^m$ are vectors whose elements are samples from a standard normal distribution.

$$y = \tilde{W}x + \tilde{b} \quad (11)$$

A high level view of a neural network with the augmented state aware perturbation module is shown in Figure 2. We partition θ into θ^b and θ^p , where θ^b is the set of parameters used to generate the hidden state representation $h(s; \theta^b)$ using the neural network h and θ^p are the parameters to which noise is to be added. Given the hidden state representation, perturbation module g_{Θ} , is used to compute the state dependent standard deviation $\sigma(h(s; \theta^b); \Theta)$, which is used to perturb the parameters θ^p of the network k . k then computes action-values for all actions. Additional features that may aid in exploration such as state visit counts or uncertainty estimates can also be appended to $h(s; \theta^b)$ before being passed as input to g_{Θ} .

Fortunato et al. [10] suggests two alternatives to generate ϵ_w and ϵ_b . The more computationally expensive alternative, Independent Gaussian noise, requires the sampling of each element of ϵ_w and ϵ_b independently, resulting in a sampling of $lm + m$ quantities per layer. Factored Gaussian noise, on the other hand, samples two vectors $\hat{\epsilon}_l$ and $\hat{\epsilon}_m$ of sizes l and m respectively. These vectors are then put through a real valued function $z(x) = \text{sgn}(x)\sqrt{x}$ before an outer product is taken to generate ϵ_w and ϵ_b (Equation 12), which are the required noise samples for the layer. Readers are referred to [10] for more details on these two noise sampling techniques. Being less computationally taxing and not having any notable impact on the performance [10], we select Factored Gaussian noise as our method

for sampling perturbations.

$$\epsilon_w = z(\hat{\epsilon}_l) \otimes z(\hat{\epsilon}_m), \quad \epsilon_b = z(\hat{\epsilon}_m) \quad (12)$$

5.2 Network Parameters and Loss Function

The set of learnable parameters for a SANE network, is a union of the set of parameters of the main network, θ , and the set of parameters of the auxiliary network perturbation module, Θ . Moreover, in place of minimizing the loss over the original set of parameters, $\mathbb{E}[\mathcal{L}(\theta)]$, the SANE network minimizes the function $\mathbb{E}[\mathcal{L}([\theta, \Theta])]$, which is the loss corresponding to the network parameterized by the perturbed weights of the network. Furthermore, with both the main network and the perturbation module being differentiable entities, using the reparameterization trick to sample the perturbations allows the joint optimization of both θ and Θ via backpropagation.

5.3 State Aware Deep Q Learning

We follow an approach similar to [10] to add state aware noisy exploration to DQNs [17]. In our implementation of a SANE DQN for Atari games, θ^b and θ^p correspond to the set of parameters in the convolutional layers and the fully connected layers respectively. The Q network and the target network have their own copies of the network and perturbation module parameters.

The DQN learns by minimizing the following loss, where θ, θ' represent the network parameters of the Q-network and the target network and Θ, Θ' represent the perturbation module parameters of the Q-network and the target network respectively. The training samples (s_i, a_i, r_i, s'_i) are drawn uniformly from the replay buffer.

$$\mathcal{L}(\theta, \Theta) = \mathbb{E} \left[\frac{1}{b} \sum_{i=1}^b (Q(s_i, a_i; \theta, \Theta) - (r_i + \gamma \max_a Q(s'_i, a; \theta', \Theta')))^2 \right]$$

5.4 Variational View of SANE DQNs

In SANE DQNs, we allow the network to use a different posterior approximation $q(\theta|s)$ for different states but restrict the perturbations that is added to all parameters to be sampled by the same distribution given a state s . Similar to NoisyNets, for the unperturbed parameters of the convolutional layers and the perturbation

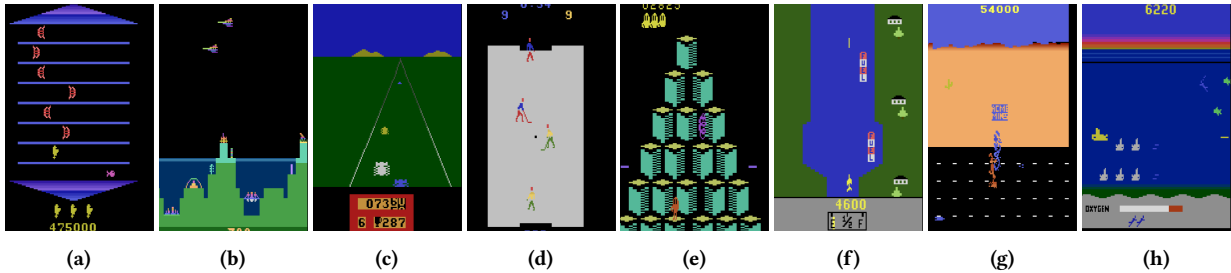


Figure 3: High risk states learnt by Q-SANE in the 8 game sub-suite

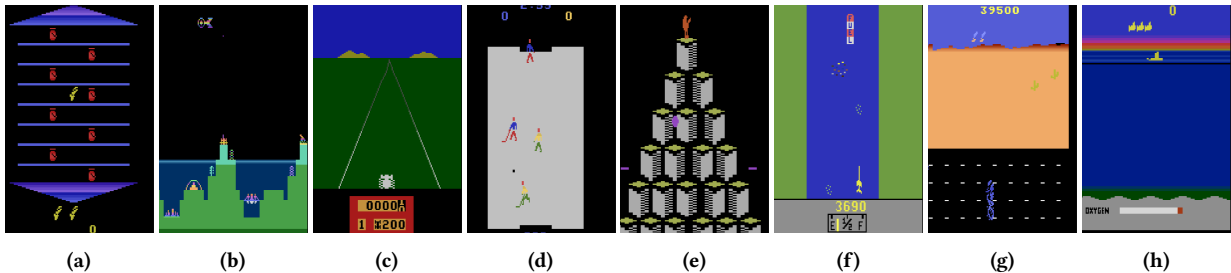


Figure 4: Low risk states learnt by Q-SANE in the 8 game sub-suite

module, we consider $q(\theta^b|s) = \mathcal{N}(\mu_b, \epsilon I)$, $q(\Theta|s) = \mathcal{N}(\mu_\Theta, \epsilon I)$ and for the parameters of the fully connected layers, we take $q(\theta^p|s) = \mathcal{N}(\mu^p, \sigma^2(h(s; \theta^b); \Theta))$. We take the prior $p(\theta) = \mathcal{N}(0, I)$ for all the parameters of the network. It follows that the objective to maximize is the same as objective (7), but where the parameters $\hat{\theta}_i$ are drawn from $q(\theta|s_i^i)$.

In our experiments, we compare two different SANE DQN variants, namely, the simple-SANE DQN and the Q-SANE DQN. Both these SANE DQNs have the same network structure as shown in Figure 2. Q-SANE DQNs and simple-SANE DQNs differ in the additional features that are added to the perturbation module. Simple-SANE DQNs add no additional features to aid the perturbation module. On the other hand, Q-SANE DQNs use the non-noisy Q-values of the state as additional features. The non-noisy Q-values are computed via a forward pass of the neural network with no perturbations applied to θ^p . Adding Q-values as features to the perturbation module can be useful, as a state where all the action values take similar values could be an indication of a low risk state and vice versa.

6 EXPERIMENTS

We conduct our experiments on a suite of 11 Atari games. This suite contains 8 games that exhibit both high and low risk states (see Figures 1, 3 and 4) that we expect would benefit from state aware exploratory behaviour. We expect SANE not to have any notable benefit in the other 3 games.

6.1 Atari Test Suite

The 11 game Atari test suite has an 8 game sub-suite, consisting of the games Asterix, Atlantis, Enduro, IceHockey, Qbert, Riverraid, RoadRunner and Seaquest. High risk and low risk states of these

games (in order) are shown in Figures 3 and 4 respectively. The games in this sub-suite have properties that benefit agents when trained with SANE exploration. Most high risk states in these games, occur when the agent is either at risk of being hit (or captured) by an enemy or at risk of going out of bounds of the play area. Figures 3a, 3b, 3c, 3e, 3g and 3h illustrate this property of high risk states. Low risk states, on the other hand, correspond to those states where the agent has a lot of freedom to move around without risking loss of life. Most of the states in Figure 4 demonstrate this.

Additionally, there may be other complex instances of high risk states. For instance, in Riverraid, states where the agent is about to run out of fuel can be considered high risk states (Figure 3f). Moreover, sometimes the riskiness of a state can be hard to identify. This is illustrated by the high and low risk states of IceHockey shown in Figures 3d and 4d respectively. In games like IceHockey, high risk states are determined by the positions of the players and the puck. Figure 3d is a high risk state as the puck’s possession is being contested by both teams, while 4d is low risk as the opponent is certain to gain the puck’s possession in the near future.

We also include 3 games, namely, FishingDerby, Boxing and Bowling, in our suite to check the sanity of SANE agents in games where we expect SANE exploration not to have any additional benefits over NoisyNets.

6.2 Parameter Initialization

We follow the initialization scheme followed by [10] to initialize the parameters of NoisyNet DQNs. Every layer of the main network of simple-SANE and Q-SANE DQNs are initialized with the Glorot uniform initialization scheme [11], while every layer of the perturbation module of both the SANE DQN variants are initialized with samples from $\mathcal{N}(0, \frac{2}{l})$, where l is number of inputs to the layer.

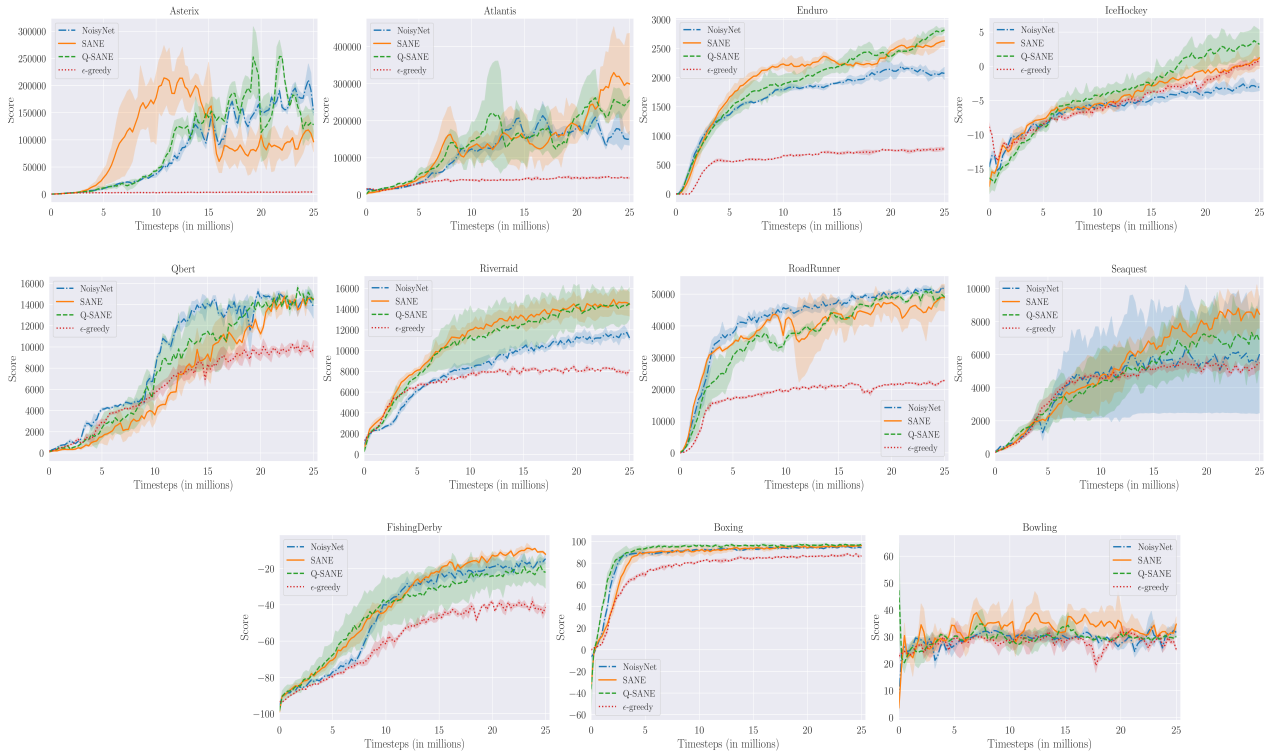


Figure 5: Learning curves of SANE DQNs, NoisyNet DQNs and ϵ -greedy DQNs.

Table 1: Scores of DQN agents when evaluated without noise injection or exploratory action selection.

Game	Q-SANE	simple-SANE	NoisyNets	ϵ -greedy
Asterix	126213±20478	133849 ± 49379	110566 ± 31800	15777 ± 3370
Atlantis	141337 ± 67719	265144 ± 151154	162738 ± 73271	229921 ± 41143
Enduro	2409±321	2798 ± 311	2075 ± 24	1736± 197
IceHockey	2.99±1.4	1.43 ± 1.76	-2.4 ± 0.52	3.46 ± 0.8
Qbert	17358 ±1015	15341 ± 162	15625 ± 166	16025 ± 555
Riverraid	14620±3491	14919 ± 997	11220 ± 223	12023 ± 512
RoadRunner	49598±1635	45929 ± 1648	51805 ± 885	47570 ± 1651
Seaquest	8368±3426	8805 ± 1392	6031 ± 3567	7682 ± 1648
FishingDerby	-19.78 ±7.2	-12.1 ± 4.2	-11.5 ± 5.4	-33.9 ± 9.1
Boxing	95.3 ± 3	93.2 ± 4.5	95.5 ± 1.7	96.6 ± 0.73
Bowling	29±1	28.08 ± 1.2	37.4 ± 3.8	20.6 ± 4.7
Score (8 games)	4.86	5.51	4.28	3.33
Score (11 games)	4.1	4.85	3.98	3.25

6.3 Architecture and Hyperparameters

We use the same network structure to train NoisyNet, simple-SANE, Q-SANE and ϵ -greedy DQNs. This network structure closely follows the architecture suggested in [17]. The inputs to all the networks are also pre-processed in a similar way.

The perturbations for NoisyNet, simple-SANE and Q-SANE DQNs are sampled using the Factored Gaussian noise setup [10]. The SANE perturbation module used for all games and all SANE

agents is a 1-hidden layer fully connected neural network. We train simple-SANE DQNs on the games of Asterix and Seaquest to determine the size of the hidden layer. A hyperparameter search over the set $\{32, 64, 128, 256, 512\}$ revealed that a module with 256 hidden neurons gave the best results on these games. The hidden layer uses ReLU activation, and output layer computes one output which corresponds to the state aware standard deviation $\sigma(h(s; \theta^b); \Theta)$.

We train the DQNs with an Adam optimizer with learning rate, $\alpha = 6.25 \times 10^{-5}$. All other hyperparameters use the same values as

Table 2: Scores of DQN Agents when evaluated with noise injection.

Game	Q-SANE	simple-SANE	NoisyNets
Asterix	182797±51182	194547 ± 56492	134682 ± 26574
Atlantis	281189±126834	230837 ± 104472	166512 ± 93945
Enduro	2849 ± 464	2855 ± 579	1946 ± 136
IceHockey	2.86 ± 1.97	1.9 ± 3.25	-1.53 ± 0.45
Qbert	16950±479	15438 ± 57	13824 ± 2690
Riverraid	15168± 2068	15434 ± 891	11076 ± 889
RoadRunner	47434± 2352	47578 ± 3787	51260 ± 712
Seaquest	7184± 2806	7844 ± 1245	6087 ± 3654
FishingDerby	-15.92 ± 7.9	-10.83 ± 2.34	-14 ± 2.8
Boxing	96.16± 1.73	95.1 ± 2.1	93.6 ± 2.7
Bowling	28.8± 0.61	28.13 ± 1.25	34.2 ± 2.4
Score (8 games)	6.43	6.2	4.64
Score (11 games)	5.54	5.37	4.22

used by [17]. The agents are trained for a total of 25M environment interactions where each training episode is limited to 100K agent-environment interactions. Both NoisyNet and SANE DQNs use greedy action selection. Please refer to Sections B and C in the Appendix for more details about the implementation.

For each game in the test suite, we train three simple-SANE, Q-SANE, NoisyNet and ϵ -greedy DQN agents. Figure 5 shows the average learning curves of all the learning agents. Each point in the learning curve corresponds to the average reward received by the agent in the last 100 episodes, averaged over 3 independent runs. Table 1 shows the mean scores achieved by simple-SANE, Q-SANE, NoisyNet and ϵ -greedy DQNs after being trained for 25M environment interactions on being evaluated for 500K frames with no noise injection. The scores of the best scoring agents in each game have been highlighted. We also evaluate simple-SANE, Q-SANE and NoisyNet DQNs with noise injection. These scores are presented in Table 2. Tables 1 and 2 also report the mean human-normalized scores (HNS) [8] achieved by these methods on the 8 games which are likely to benefit from SANE exploration and on the whole 11 game suite. We also present some high-risk and low-risk states identified by Q-SANE agents in Figures 3 and 4.

We observe that when evaluated without noise injection, both Q-SANE and simple-SANE outperform NoisyNets in 6 of the 8 games in the sub-suite. NoisyNets achieve higher scores than both SANE variants in RoadRunner. In the three games not in the sub-suite, NoisyNets achieve higher but similar scores in FishingDerby and Boxing while performing much better in Bowling compared to the other agents. Evaluating the agents with noise injection proves beneficial for both SANE and NoisyNet agents, all of them achieving higher mean HNS in the 8 game sub-suite and the whole test suite. However, simple-SANE and Q-SANE agents achieve greater gains as they score higher than NoisyNets in 7 games in the sub-suite. SANE agents also score better in the remaining three games but do not manage to score better than NoisyNets in Bowling. Q-SANE and simple-SANE achieve the highest mean HNS on both the 8 game sub-suite and the whole test suite when evaluated with and without noise injection respectively.

7 CONCLUSION

In this paper, we derive a variational Thompson sampling approximation for DQNs, which uses the distribution over the network parameters as a posterior approximation. We interpret NoisyNet DQNs as an approximation to this variational Thompson Sampling method where the posterior is approximated using a state uniform Gaussian distribution. Using a more general posterior approximation, we propose State Aware Noisy Exploration, a novel exploration strategy that enables state dependent exploration in deep reinforcement learning algorithms by perturbing the model parameters. The perturbations are injected into the network with the help of an augmented SANE module, which draws noise samples from a Gaussian distribution whose variance is conditioned on the current state of the agent. We hypothesize that such state aware perturbations are useful to direct exploration in tasks that go through a combination of high risk and low risk situations, an issue not considered by other methods that rely on noise injection.

We test this hypothesis by evaluating two SANE DQN variants, namely simple-SANE and Q-SANE DQNs, on a suite of 11 Atari games containing a selection of games, most of which fit the above criterion and some that do not. We observe that both simple-SANE and Q-SANE perform better than NoisyNet agents in most of the games in the suite, achieving better mean human-normalized scores.

An additional benefit of SANE noise injection mechanism is its flexibility of design. SANE effects exploration via a separate perturbation module, the size or architecture of which is not tied to the model being perturbed and hence is flexible to user design and can be tailored to the task. As a consequence, this exploration method might scale better to larger network models. Hence, SANE presents a computationally inexpensive way to incorporate state information into exploration strategies and is a step towards more effective, efficient and scalable exploration.

ACKNOWLEDGEMENTS

This work was supported by the National Research Foundation Singapore under its AI Singapore Program (Award Number: AISGRP-2018-006).

REFERENCES

- [1] Joshua Achiam and Shankar Sastry. 2017. Surprise-Based Intrinsic Motivation for Deep Reinforcement Learning. *arXiv preprint arXiv:1703.01732* (2017).
- [2] Shipra Agrawal and Navin Goyal. 2012. Analysis of Thompson Sampling for the Multi-Armed Bandit Problem. In *Conference on Learning Theory*. 39–1.
- [3] Shipra Agrawal and Navin Goyal. 2013. Further Optimal Regret Bounds for Thompson Sampling. In *Artificial Intelligence and Statistics*. 99–107.
- [4] Shipra Agrawal and Randy Jia. 2017. Optimistic Posterior Sampling for Reinforcement Learning: Worst-Case Regret Bounds. In *Advances in Neural Information Processing Systems*. 1184–1194.
- [5] Peter Auer and Ronald Ortner. 2007. Logarithmic Online Regret Bounds for Undiscounted Reinforcement Learning. In *Advances in Neural Information Processing Systems*. 49–56.
- [6] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. 2018. Efficient Exploration through Bayesian Deep Q-Networks. In *2018 Information Theory and Applications Workshop (ITA)*. IEEE, 1–9.
- [7] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. 2016. Unifying Count-Based Exploration and Intrinsic Motivation. In *Advances in Neural Information Processing Systems*. 1471–1479.
- [8] Marc Brittain, Joshua R. Bertram, Xuxi Yang, and Peng Wei. 2019. Prioritized Sequence Experience Replay. *CoRR abs/1905.12726* (2019). [arXiv:1905.12726](http://arxiv.org/abs/1905.12726) <http://arxiv.org/abs/1905.12726>
- [9] Carlos Florensa, Yan Duan, and Pieter Abbeel. 2017. Stochastic Neural Networks for Hierarchical Reinforcement Learning. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=B1oK8aoxe>
- [10] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Matteo Hessel, Ian Osband, Alex Graves, Volodymyr Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, Charles Blundell, and Shane Legg. 2018. Noisy Networks For Exploration. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rywHCPkAW>
- [11] Xavier Glorot and Yoshua Bengio. 2010. Understanding the Difficulty of Training Deep Feedforward Neural Networks. In *Proceedings of The Thirteenth International Conference on Artificial Intelligence and Statistics*. 249–256.
- [12] Shixiang Gu, Ethan Holly, Timothy Lillicrap, and Sergey Levine. 2017. Deep Reinforcement Learning for Robotic Manipulation with Asynchronous Off-Policy Updates. In *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 3389–3396.
- [13] Thomas Jaksch, Ronald Ortner, and Peter Auer. 2010. Near-Optimal Regret Bounds for Reinforcement Learning. *Journal of Machine Learning Research* 11, Apr (2010), 1563–1600.
- [14] Emilie Kaufmann, Nathaniel Korda, and Rémi Munos. 2012. Thompson Sampling: An Asymptotically Optimal Finite-Time Analysis. In *International Conference on Algorithmic Learning Theory*. Springer, 199–213.
- [15] Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. 2017. Learning Combinatorial Optimization Algorithms over Graphs. In *Advances in Neural Information Processing Systems*. 6348–6358.
- [16] Xiaodan Liang, Lisa Lee, and Eric P Xing. 2017. Deep Variation-Structured Reinforcement Learning for Visual Relationship and Attribute Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 848–857.
- [17] 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. 2015. Human-Level Control through Deep Reinforcement Learning. *Nature* 518, 7540 (2015), 529.
- [18] Ian Osband, John Aslanides, and Albin Cassirer. 2018. Randomized Prior Functions for Deep Reinforcement Learning. In *Advances in Neural Information Processing Systems*. 8617–8629.
- [19] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. 2016. Deep Exploration via Bootstrapped DQN. In *Advances in Neural Information Processing Systems*. 4026–4034.
- [20] Ian Osband and Benjamin Van Roy. 2015. Bootstrapped Thompson Sampling and Deep Exploration. *arXiv preprint arXiv:1507.00300* (2015).
- [21] Ian Osband and Benjamin Van Roy. 2017. Why is Posterior Sampling Better than Optimism for Reinforcement Learning?. In *International Conference on Machine Learning*. 2701–2710.
- [22] Ian Osband, Benjamin Van Roy, and Zheng Wen. 2016. Generalization and Exploration via Randomized Value Functions. In *Proceedings of the 33rd International Conference on Machine Learning-Volume 48*. JMLR. org, 2377–2386.
- [23] Georg Ostrovski, Marc G Bellemare, Aaron van den Oord, and Rémi Munos. 2017. Count-Based Exploration with Neural Density Models. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2721–2730.
- [24] Matthias Plappert, Rein Houthoofd, Prafulla Dhariwal, Szymon Sidor, Richard Y. Chen, Xi Chen, Tamim Asfour, Pieter Abbeel, and Marcin Andrychowicz. 2018. Parameter Space Noise for Exploration. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=ByBA12eAZ>
- [25] Thomas Rückstieß, Frank Sehne, Tom Schaul, Daan Wierstra, Yi Sun, and Jürgen Schmidhuber. 2010. Exploring Parameter Space in Reinforcement Learning. *Paladyn* 1, 1 (2010), 14–24.
- [26] Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor, and Ilya Sutskever. 2017. Evolution Strategies as a Scalable Alternative to Reinforcement Learning. *arXiv preprint arXiv:1703.03864* (2017).
- [27] Aravindan Siddharth. 2021. SANE. <https://github.com/NUS-Learning-Inference-Decision-Group/SANE>.
- [28] Bradley C Stadie, Sergey Levine, and Pieter Abbeel. 2015. Incentivizing Exploration in Reinforcement Learning with Deep Predictive Models. *arXiv preprint arXiv:1507.00814* (2015).
- [29] Malcolm Strens. 2000. A Bayesian Framework for Reinforcement Learning. In *International Conference on Machine Learning*, Vol. 2000. 943–950.
- [30] Richard S Sutton and Andrew G Barto. 2018. *Reinforcement learning: An Introduction*. MIT press.
- [31] Adrien Ali Taiga, William Fedus, Marlos C. Machado, Aaron Courville, and Marc G. Bellemare. 2020. On Bonus Based Exploration Methods In The Arcade Learning Environment. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=BJewlyStDr>
- [32] William R Thompson. 1933. On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples. *Biometrika* 25, 3/4 (1933), 285–294.
- [33] Sirui Xie, Junjing Huang, Lanxin Lei, Chunxiao Liu, Zheng Ma, Wei Zhang, and Liang Lin. 2019. NADPEX: An On-Policy Temporally Consistent Exploration Method for Deep Reinforcement Learning. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rkxciiC9tm>

Algorithm 1 SANE Deep Q Learning

```
1: Initialize the target network parameters  $\hat{\theta}, \hat{\Theta}$  and the Q network
   parameters  $\theta, \Theta$  randomly.
2: Initialize an empty experience replay buffer;
3: Initialize noise sampling method  $N$ .
4: steps  $\leftarrow 0$ 
5: while steps  $\leq$  max_steps do
6:    $t \leftarrow 0$ 
7:   Observe  $s_0$ 
8:   while  $s_t$  not terminal do
9:     Compute  $h(s_t; \theta^b)$ .
10:    Sample perturbations  $\epsilon \sim N(0, 1)$  to perturb  $\theta^P$ 
11:     $\theta^{P'} \leftarrow \theta^P + \sigma(h(s_t; \theta^b); \Theta)\epsilon$ 
12:     $a_t \leftarrow \max_a Q(s_t, a; \theta^{P'}, \theta^b)$ 
13:    Take  $a_t$ , observe next state  $s'$ , reward  $r_t$ 
14:    Add transition  $(s_t, a_t, r_t, s')$  to the replay buffer;
15:    if buffer_size  $\geq$  max_buffer_size then
16:      Remove oldest buffer entry;
17:    if steps mod update_frequency = 0 then
18:      Sample  $b$  transitions  $\{(s_i, a_i, r_i, s'_i), 1 \leq i \leq b\}$  uni-
   formly from the replay buffer.
19:       $L \leftarrow 0$ 
20:      for  $i \in \{1 \dots b\}$  do
21:        Compute  $h(s'_i; \theta^b)$  and  $h(s'_i; \hat{\theta}^b)$ 
22:        Sample  $\epsilon, \hat{\epsilon} \sim N(0, 1)$  to perturb  $\theta^P, \hat{\theta}^P$ 
23:         $\theta^{P'} \leftarrow \theta^P + \sigma(h(s_t; \theta^b); \Theta)\epsilon$ 
24:         $\hat{\theta}^{P'} \leftarrow \hat{\theta}^P + \sigma(h(s_t; \hat{\theta}^b); \hat{\Theta})\hat{\epsilon}$ 
25:         $T = (r_i + \gamma \max_a Q(s'_i, a; \theta^{P'}, \hat{\theta}^b))$ 
26:         $L = L + \frac{1}{b} \left[ Q(s_i, a_i; \theta^{P'}, \theta^b) - T \right]^2$ 
27:      Backpropagate to minimize the batch loss  $L$ 
28:      if steps mod copy_frequency = 0 then
29:         $\hat{\theta} \leftarrow \theta; \hat{\Theta} \leftarrow \Theta$ 
30:       $t \leftarrow t + 1$ ; steps  $\leftarrow$  steps+1
```

A DERIVING THE ELBO

In Section 4, we mentioned that minimizing the $KL(q(\theta), p(\theta|D))$ is equivalent to maximizing the ELBO. Here we provide a derivation of that claim.

$$KL(q(\theta), p(\theta|D)) = \int q(\theta) \log \frac{q(\theta)}{p(\theta|D)} d\theta$$

Now,

$$\begin{aligned} p(\theta|D) &= p(\theta|X, Y) = \frac{p(\theta)p(X, Y|\theta)}{p(X, Y)} = \frac{p(\theta)p(Y|X, \theta)p(X|\theta)}{p(X, Y)} \\ &= \frac{p(\theta)p(Y|X, \theta)p(X)}{p(X, Y)}; \end{aligned}$$

Substituting the above value,

$$\begin{aligned} KL(q(\theta), p(\theta|D)) &= \int q(\theta) \left[\log \frac{q(\theta)p(X, Y)}{p(\theta)p(Y|X, \theta)p(X)} \right] d\theta \\ &= \int q(\theta) \log \frac{q(\theta)}{p(\theta)p(Y|X, \theta)} d\theta + \int q(\theta) \log \frac{P(X, Y)}{P(X)} d\theta \\ &= \int q(\theta) \log \frac{q(\theta)}{p(\theta)p(Y|X, \theta)} d\theta + \log \frac{P(X, Y)}{P(X)} \\ &= \int q(\theta) \log \frac{q(\theta)}{p(\theta)} d\theta - \int q(\theta) \log p(Y|X, \theta) d\theta + \log \frac{P(X, Y)}{P(X)} \\ &= KL(q(\theta), p(\theta)) - \int q(\theta) \log p(Y|X, \theta) d\theta + \log \frac{P(X, Y)}{P(X)} \end{aligned}$$

Now, $P(X, Y)$ and $P(X)$ are fixed with respect to θ . So, minimizing $KL(q(\theta), p(\theta|D))$ is equivalent to maximizing the ELBO (Equation 2).

B THE SANE-DQN ALGORITHM

Algorithm 1 describes the implementation of a SANE DQN. The Q network and the target network in SANE DQNs have their own copies of the network parameters. These are denoted by θ and $\hat{\theta}$ respectively. The Q and target networks also maintain different copies of the parameters of the perturbation module, Θ and $\hat{\Theta}$ respectively. Further, θ is partitioned into two sets, θ^b, θ^P , where θ^b is the set of *base* parameters that help us compute the hidden state representation $h(s; \theta^b)$ and θ^P is the set of network parameters that are to be perturbed with state aware perturbations. We define similar counterparts $\hat{\theta}^b$ and $\hat{\theta}^P$ in the target network.

With every forward pass, the network first calculates $h(s, \theta^b)$, which is then passed to the perturbation module. Factored Gaussian noise samples are procured and multiplied with $\sigma(h(s; \theta^b); \Theta)$, to get perturbations equivalent to those directly sampled from $\mathcal{N}(0, \sigma^2(h(s; \theta^b); \Theta))$ (Equation 10). These perturbations are added to the parameters θ^P . The agent then selects the action greedily with respect to the action values computed by the perturbed Q-network. While computing the batch-loss, the Q network and target network parameters θ^P and $\hat{\theta}^P$ are perturbed with state aware perturbations sampled from $\mathcal{N}(0, \sigma^2(h(s; \theta^b); \Theta))$ and $\mathcal{N}(0, \sigma^2(h(s; \hat{\theta}^b); \hat{\Theta}))$ respectively (Lines 21-25 in Algorithm 1). This loss is then backpropagated to train θ and Θ .

C DQN IMPLEMENTATION DETAILS

The network structure and input pre-processing closely follows the architecture and method suggested in [17]. The DQN consists of three convolutional layers followed by two linear layers. The first convolutional layer has 32 filters of size 8×8 . This layer is followed by a convolutional layer with 64 filters of size 4×4 . The last convolutional layer has 64 filters of size 3×3 . The convolutional layers use strides of 4, 2 and 1 respectively. The convolutional layers are followed by 2 fully connected layers, a hidden layer with 512 neurons and an output layer with the number of outputs being equal to the number of actions available to the agent for the task. With the exception of the output layer, a ReLU activation function follows every layer.

The perturbations for both the networks are sampled using the Factored Gaussian noise setup [10]. The state aware perturbation module used for all games is a 1-hidden layer fully connected neural network. The hidden layer consists of 256 neurons and uses ReLU activation. The output layer computes one output which corresponds to the state aware standard deviation $\sigma(h(s; \theta^b); \Theta)$.

We train the DQNs with an Adam optimizer with learning rate, $\alpha = 6.25 \times 10^{-5}$ and $\epsilon = 1.5 \times 10^{-4}$. We use a replay buffer that can hold a maximum of 1M transitions. We populate the replay buffer with 50K transitions that we obtain by performing random actions for ϵ -greedy agents and by following the policy suggested by the network for NoisyNet and SANE agents. Thereafter, we train the Q network once every 4 actions, with a batch of 32 transitions sampled uniformly from the replay buffer. We copy over the parameters of the Q network to the target network after every 10K transitions. We use a discount factor of $\gamma = 0.99$ for all games. Additionally, the rewards received by the agent are clipped in the range $[-1, 1]$.

The agent is trained for a total of 25M agent-environment interactions. The input to the network is a concatenation of 4 consecutive frames, and we take a random number of no-op actions (upto 30) at the start of each episode, so that the agent is given a random start.

The codebase for our SANE, Q-SANE, NoisyNet and ϵ -greedy DQN agents are available at [27].

D ADDITIONAL EXPERIMENTAL DETAILS

D.1 Human Normalized Scores

The human normalized score for any agent is calculated as follows

$$\text{HNS} = \frac{\text{Score}_{\text{agent}} - \text{Score}_{\text{random}}}{\text{Score}_{\text{human}} - \text{Score}_{\text{random}}} \quad (13)$$

We use the same random and human baseline scores as used in [8]. We list these baseline scores in Table 3 for easy access.

Table 3: Baseline human and random values used to calculate Human Normalized Scores

Game	Human Score	Random Score
Asterix	8503	210
Atlantis	29028	12850
Boxing	12.1	0.1
Bowling	160.7	23.1
Enduro	860.5	0
FishingDerby	-38.7	-91.7
IceHockey	0.9	-11.2
Qbert	13455	163.9
Riverraid	17118	1338.5
RoadRunner	7845	11.5
Seaquest	42054	68.4

D.2 Visualizations

We present the high risk and low risk states (from Figures 3 and 4) along with the state aware standard deviation predicted by the Q-SANE agents in Figures 6 and 7 respectively.

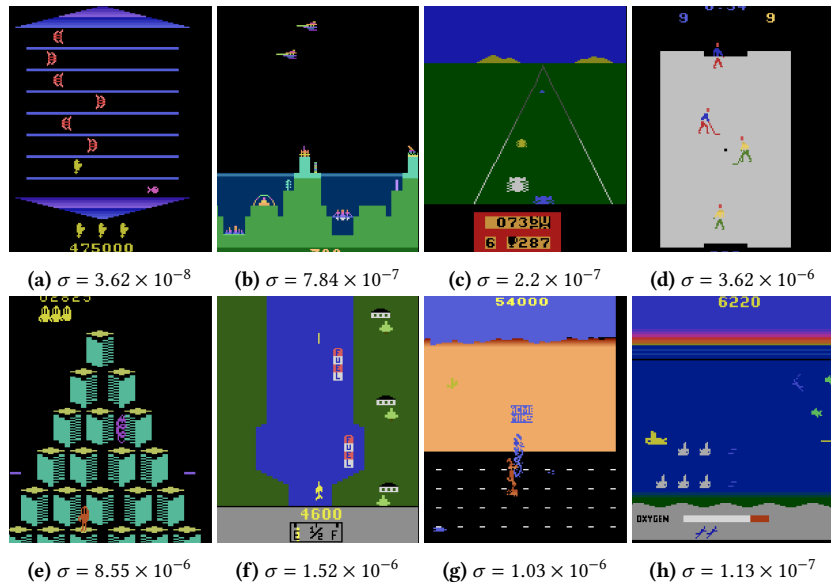


Figure 6: High risk states learnt by Q-SANE in the 8 game sub-suite. The captions mention the standard deviation predicted for each state.

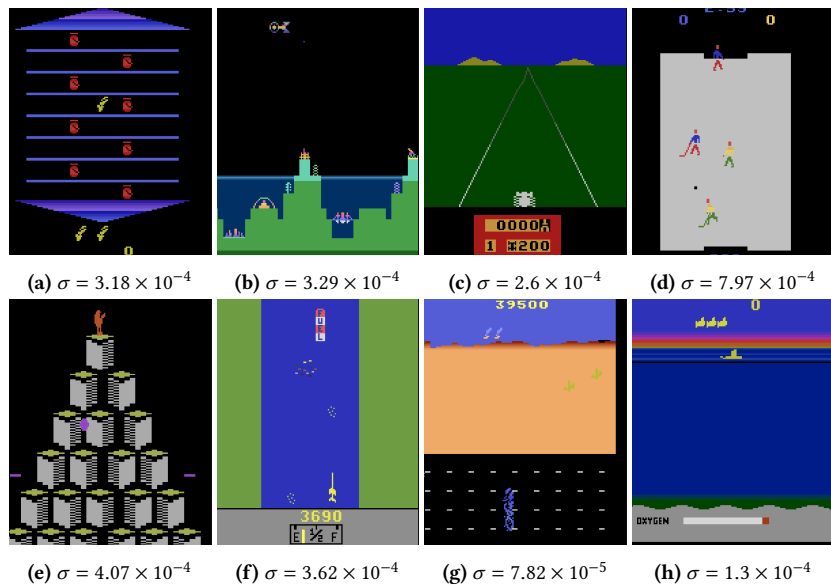


Figure 7: Low risk states learnt by Q-SANE in the 8 game sub-suite. The captions mention the standard deviation predicted for each state.