

# VASE: Variational Assorted Surprise Exploration for Reinforcement Learning

Haitao Xu, Brendan McCane, Lech Szymanski  
 Department of Computer Science  
 University of Otago  
 Dunedin, New Zealand  
 {haitao,mccane,lechsyzm}@cs.otago.ac.nz

arXiv:1910.14351v1 [cs.LG] 31 Oct 2019

**Abstract**—Exploration in environments with continuous control and sparse rewards remains a key challenge in reinforcement learning (RL). Recently, surprise has been used as an intrinsic reward that encourages systematic and efficient exploration. We introduce a new definition of surprise and its RL implementation named Variational Assorted Surprise Exploration (VASE). VASE uses a Bayesian neural network as a model of the environment dynamics and is trained using variational inference, alternately updating the accuracy of the agent’s model and policy. Our experiments show that in continuous control sparse reward environments VASE outperforms other surprise-based exploration techniques.

**Index Terms**—surprise, reinforcement learning, exploration, variational inference, Bayesian neural network

## I. INTRODUCTION

Reinforcement learning (RL) trains agents to act in an environment so as to maximise cumulative reward. The resulting behaviour is highly dependent on the trade-off between exploration and exploitation. During training, the more the agent departs from its current policy, the more it learns about the environment, which may lead it to a better policy; the closer it adheres to the current policy, the less time wasted exploring less effective options. How much and where to explore has an immense impact on the training and ultimately on what the agent learns. Designing exploration strategies, especially for increasingly complex environments, is still a significant challenge.

A common approach to exploration strategies is to rely on heuristics that introduce random perturbations into the choices of actions during training, such as  $\epsilon$ -greedy [1] or Boltzmann exploration [2]. These methods instruct the agent to occasionally take an arbitrary action that may drive it into a new experience. Another way is through the addition of noise to the parameter space of the agent’s policy neural network [3, 4], which varies the policy itself to a similar random exploration net effect. These strategies can be highly inefficient because they are a result of random behaviour, which is especially problematic in high dimension state-action spaces (common in discretised continuous state-action space environments) because of the curse of dimensionality. Random exploration is also extremely inefficient in environments with sparse rewards, where the agent ends up wandering aimlessly through the state-space, learning nothing until (by sheer luck)

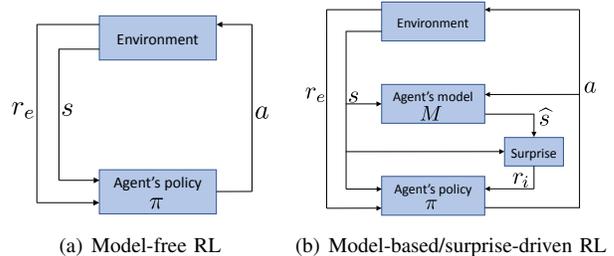


Fig. 1. Model-free RL vs. model-based/surprise-driven RL with  $s$  for state,  $r_e$  for the extrinsic (environment-driven) reward,  $a$  for action,  $r_i$  for intrinsic (agent-driven) reward,  $\pi$  for the policy, and  $M$  for the model of the environment that makes a prediction of the next state  $\hat{s}$ .

it chances upon a reward. Not surprisingly, more methodical approaches were devised, which provide the agent with intrinsic rewards that encourage efficient exploration. These intrinsic rewards are derived from computations related to the notion of curiosity [5, 6] or surprise [7, 8].

In this paper, we propose a new definition of surprise, which drives our agents’ intrinsic reward function. To compute and use this surprise for guiding exploration, we propose an algorithm called VASE (Variational Assorted Surprise Exploration) in a model-based RL framework (see Figure 1). VASE alternates the update step between the agent’s policy and its model of the environment dynamics. The policy is implemented with a multilayer feed-forward (MLFF) neural network and the dynamics model with a Bayesian neural network (BNN [9, 10, 11]). We evaluate the performance of our method against other surprise-driven methods on continuous control tasks with sparse rewards. Our experimental results show VASE’s superior performance.

## II. RELATED WORK

In RL a finite horizon discounted Markov decision process (MDP) is defined by a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma, T, \rho_0)$  where:  $\mathcal{S}$  is a state set,  $\mathcal{A}$  an action set,  $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is a transition probability distribution,  $r : \mathcal{S} \rightarrow \mathbb{R}$  a reward function,  $\gamma \in (0, 1]$  a discount factor,  $T$  the horizon and  $\rho_0$  an initial state distribution. A policy  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  gives the probability with  $\pi(a|s)$  of taking action  $a$  in state  $s$ . Let  $\tau = (s_0, a_0, \dots)$  denote the whole trajectory,  $s_0 \sim \rho_0$ ,

$a_t \sim \pi(a_t|s_t)$ . Our purpose is to find a policy  $\pi$ , modelled by  $f_\pi(s_t, \psi)$  with parameters  $\psi$ , which maximises the expected discounted total return  $\mathbb{E}_\tau[\sum_{t=0}^T \gamma^t r(s_t)]$  – a discounted sum of all rewards in a fixed horizon  $T$ .

### A. Intrinsic rewards

Intrinsic motivation is essential for effective exploration when training the agent in an environment with sparse extrinsic rewards, or no rewards at all. The overall reward signal is computed as follows:

$$r(s_t) = r_e(s_t) + r_i(s_t),$$

where  $r_e(s_t)$  represents the extrinsic reward from the environment,  $r_i(s_t)$  represents the intrinsic reward computed by the agent. Even when  $r_e(s_t) = 0$ , the intrinsic reward contributes to the cumulative reward, thus driving the learning process until non-zero extrinsic rewards are found. There are two broad approaches for encouraging the agent to explore through intrinsic rewards.

The first is the count-based approach [12, 13, 14, 15], which maintains visit counters over all the states. The intrinsic reward is inversely proportional to the current state’s counter, thus rewarding exploration of less frequented states. This approach becomes intractable with the increase of possible states, making it unfit for scenarios that have continuous state-action spaces. The second is surprise-based.

### B. Surprise-driven reinforcement learning

In model-based RL the agent maintains one or more models that predict the next state based on the current state and the action about to be taken. There are four possible regimes involving increasing amounts of probabilistic reasoning:

- 1) There is a single deterministic model
- 2) There is one model that produces a distribution over new states
- 3) There is a distribution of models, each of which is deterministic
- 4) There is a distribution of models, each of which produces a distribution over states.

As an example, the model for case 1 could be a traditional neural network, case 2 variational auto-encoder, case 3 Bayesian neural network and case 4 Bayesian variational auto-encoder. In cases 2-3, the outcome is a distribution over states and we are free to choose the most convenient formalism. For this paper we choose case 3 as in [16]. In this case the agent maintains a distribution  $P(M)$  over models or hypotheses  $M \in \mathcal{M}$ , where  $M : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$  predicts state  $s_{t+1}$ , given state  $s_t$  and action  $a_t$ . Furthermore, we have:

$$P(s_{t+1}|s_t, a_t, M) = \begin{cases} P(M|s_t, a_t), & \text{when } s_{t+1} = M(s_t, a_t) \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

and

$$P(s_{t+1}|s_t, a_t) = \int_{\mathcal{M}} P(M|s_t, a_t)P(s_{t+1}|s_t, a_t, M)\partial M. \quad (2)$$

The environment is assumed to be stochastic, but we consider the distributions in Equations 1 and 2 are subjective agent-based beliefs about the environment, and not the underlying objective truth. Furthermore, these distributions are non-stationary and change as the agent learns more about the environment. Also note that we equate states with observations even in complicated scenarios when the true state of the environment is not directly observed (e.g. the agent has a camera). This is because the goal of the agent in RL is to maximise long-term rewards, and we assume that those rewards are directly observable from the environment.

The distribution in Equation 2 allows for a number of definitions of surprise, some of which have previously appeared in the literature. In all cases, the accuracy of the distribution, after observing the next state, is used to derive surprise,  $U(s_{t+1})$ , and forms the basis for intrinsic reward:  $r_i(s_t) = \eta U(s_t)$ , where  $0 < \eta \leq 1$ , and is used to encourage exploration of *unexpected* states [7, 8, 5, 17, 18, 19]. The agent is rewarded for curiosity of the *unknown* as gauged by its model of the environment. Throughout training, the model is improved to be more accurate (and so less surprised) next time it encounters an already explored state. The hope is that this methodical approach to exploration will result in a speedier arrival of the agent at the states with non-zero extrinsic reward.

It is not entirely obvious how best to define the surprise. One proposed definition is the so-called *surprisal* [20], which is the negative log-likelihood (NLL) of the next state in RL tasks

$$U_{\text{NLL}}(s_{t+1}^*) = -\log P(s_{t+1} = s_{t+1}^*|s_t, a_t), \quad (3)$$

where  $s_{t+1}^*$  is the observed state at time  $t + 1$ . This notation is a bit clumsy, so we shorten it to:

$$U_{\text{NLL}}(s_{t+1}) = -\log P(s_{t+1}|s_t, a_t), \quad (4)$$

which we think is clear and more concise.

Surprisal is intuitive, simple and easy to compute, but it does not capture all the information available, because it only measures the surprise at a single point. An alternative is *Bayesian surprise*, which measures the difference between the prior distribution over the model space at time  $t$ , and the posterior distribution updated by Bayes’ rule with the newly observed state (first used in RL by Storck et al. [21]):

$$P(M|s_t, a_t, s_{t+1}) = \frac{P(M|s_t, a_t)P(s_{t+1}|s_t, a_t, M)}{P(s_{t+1}|s_t, a_t)} \quad (5)$$

Bayesian surprise is defined as the Kullback-Leiber (KL) divergence between the prior and the posterior beliefs about the dynamics of the environment [22, 23]:

$$U_{\text{Bayes}}(s_{t+1}) = D_{\text{KL}}[P(M|s_t, a_t)||P(M|s_t, a_t, s_{t+1})]. \quad (6)$$

Bayesian surprise measures the difference between subjective beliefs prior and post an event. One problem with Bayesian surprise is that the agent does not express surprise until it updates its belief, which is inconsistent with the instantaneous response to surprise displayed by neural data [24].

Faraji et al. [25] introduced a modification of Bayesian surprise referred to as the *confidence-corrected surprise* (CC), which measures the difference between the agent’s current beliefs about the world, and a naive observer who believes all models are equally likely:

$$U_{CC}(s_{t+1}) = D_{KL}[P(M|s_t, a_t) || P^{\text{flat}}(M|s_{t+1})], \quad (7)$$

They also designed a surprise minimization rule to let the agents adapt quickly to the environment, especially in a highly volatile environment. But confidence-corrected surprise was not applied to RL frameworks. Nevertheless, their definition of surprise is similar in spirit to ours.

To use Bayesian surprise for RL tasks, Houthoof et al. [7] proposed a surprise-driven exploration strategy called VIME (variational information maximizing exploration). They showed that VIME achieves significantly better performance compared to heuristic exploration methods across a variety of continuous control tasks with sparse rewards. However, to compute each reward, VIME needs to calculate the gradient through a Bayesian neural network (BNN)[9, 10], which is used to implement the agent’s model. This requires a forward and a backward pass, which leads to slow training speed. Achiam and Sastry [8] chose the surprisal of new observation as their intrinsic surprise reward. Their experimental results showed that surprisal does not perform as well as VIME, but it runs faster. They also showed that the surprisal includes  $L_2$ -squared model prediction error, which was first proposed in [26] and later used as the curiosity in [5]. Burda et al. [6] also chose surprisal-based strategy exploration in large scale RL environments that provide no extrinsic rewards.

### III. ASSORTED SURPRISE FOR EXPLORATION

In this paper, we focus on RL environments with continuous control and very sparse extrinsic rewards. To conquer the numerous shortcomings of existing definitions for surprise discussed in Section II-B, and inspired by the idea of confidence-corrected surprise [25], we note that in RL the definition of surprise should have the following characteristics:

- 1) subjectivity – the agent should hold subjective beliefs about the environment captured through  $P(M)$ ; the surprise depends on an agent’s belief and this belief can be updated during learning;
- 2) consistency – based on its belief, the agent should be more surprised by states with lower likelihood;
- 3) instancy – the agent should be surprised immediately when it observes a new state from the environment, without the need to update its belief first.

In order to address all of the above characteristics, we propose a new definition of surprise, which we refer to as *assorted surprise for exploration* (ASE):

$$U_{ASE}(s_{t+1}) = \mathbb{E}_{M \sim P(\cdot|s_t, a_t)}[-\log P(s_{t+1}|s_t, a_t, M)] - \delta H(P(M|s_t, a_t)), \quad (8)$$

where  $\delta$  is a trade-off coefficient. The first term,  $\mathbb{E}_{M \sim P(\cdot|s_t, a_t)}[-\log P(s_{t+1}|s_t, a_t, M)]$ , we call the assorted

surprise term, and the second term,  $H(P(M|s_t, a_t))$ , the confidence term.

We will next demonstrate that the assorted surprise satisfies all three characteristics we mentioned above. Subjectivity comes from the assorted surprise term in Equation 8 because it is an expectation over the agent’s belief in the veracity of each of the models. It is also the sum of the Bayesian surprise and the surprisal (hence the name *assorted*) as shown in the following lemma.

**Lemma 1** (Assorted surprise is the sum of Bayesian surprise and surprisal).

$$\mathbb{E}_{M \sim P(\cdot|s_t, a_t)}[-\log P(s_{t+1}|s_t, a_t, M)] = U_{Bayes}(s_{t+1}) + U_{NLL}(s_{t+1})$$

*Proof.*

$$\begin{aligned} & \mathbb{E}_{M \sim P(\cdot|s_t, a_t)}[-\log P(s_{t+1}|s_t, a_t, M)] \\ &= - \int_{\mathcal{M}} P(M|s_t, a_t) \log P(s_{t+1}|s_t, a_t, M) \partial M \\ &= - \int_{\mathcal{M}} P(M|s_t, a_t) \log \frac{P(M|s_t, a_t, s_{t+1})P(s_{t+1}|s_t, a_t)}{P(M|s_t, a_t)} \partial M \\ &= \int_{\mathcal{M}} P(M|s_t, a_t) \log \frac{P(M|s_t, a_t)}{P(M|s_t, a_t, s_{t+1})} \partial M \\ &\quad - \int_{\mathcal{M}} P(M|s_t, a_t) \log P(s_{t+1}|s_t, a_t) \partial M \\ &= D_{KL}[P(M|s_t, a_t) || P(M|s_t, a_t, s_{t+1})] \\ &\quad - \log P(s_{t+1}|s_t, a_t) \\ &= U_{Bayes}(s_{t+1}) \\ &\quad + U_{NLL}(s_{t+1}), \end{aligned}$$

□

This means that assorted surprise is subjective due to the contribution of  $U_{Bayes}$ . However, the expectation in Eq. 8 does not require evaluation of  $P(M|s_t, a_t, s_{t+1})$ , so there is no requirement to update the agent’s belief in order to compute assorted surprise thus satisfying the instancy characteristic. Finally, a less likely state leads to a larger surprise through the negative log likelihood of  $P(s_{t+1}|s_t, a_t)$  in the surprisal term. This satisfies the consistency requirement.

The confidence term of Eq. 8,  $H(P(M|s_t, a_t))$ , is the entropy of  $P(M|s_t, a_t)$ . This term was added for confidence correction of assorted surprise. A confident agent will have a low entropy, and therefore any surprising event according to the assorted surprise term will remain surprising. Whereas an uncertain agent will have a large entropy and their overall surprise will be reduced because they would be equally surprised by many events. That is, confident agents are more surprised when their beliefs are violated by unlikely events than uncertain agents.

Assorted surprise captures the positive elements of both Bayesian surprise and confidence-corrected surprise. It implicitly computes the difference in belief as in Bayesian surprise without needing to update the belief first, and it can be

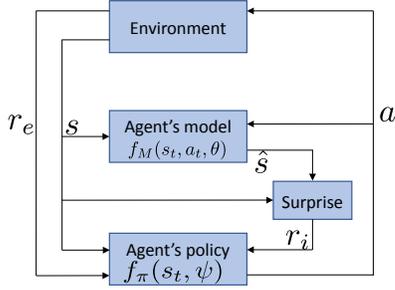


Fig. 2. VASE-driven RL with  $s$  for state,  $r_e$  for the extrinsic (environment-driven) reward,  $a$  for action,  $r_i$  for intrinsic (agent-driven) reward,  $a_t = f_\pi(s_t, \psi)$  the policy governed by set of parameters  $\psi$ , and  $\hat{s} = f_m(s_t, a_t, \theta)$  the sample BNN model of the environment.

computed very fast as in confidence-corrected surprise without needing to maintain the idea of a naive observer.

### A. Variational Assorted Surprise for Exploration (VASE)

Based on the discussion in Section II-B, we construct a BNN dynamics model  $f_m(s_t, a_t, \Theta)$ , where  $\Theta$  is a random variable describing the parameters of the model (see Figure 2). BNN can be seen as a distribution of models  $M$ , where a sample of network parameters  $\theta$  according to distribution  $P(\theta)$  is analogous to generating a single prediction of the next state according to  $P(M)$ . The prior distribution  $P(\theta)$  changes to posterior  $P(\theta|\mathcal{D})$  when BNN is trained by  $\mathcal{D} = \{s_t, a_t, s_{t+1}\}$ .

Since the posterior  $P(\theta|s_t, a_t)$  in Eq. 8 is intractable, we turn to variational inference [27] to approximate it with a fully factorised Gaussian distribution [9, 10, 11]

$$q(\theta; \phi) = \prod_{i=1}^{|\Theta|} \mathcal{N}(\theta_i; \mu_i, \sigma_i^2), \quad (9)$$

where  $\theta_i$  is the  $i^{\text{th}}$  component of  $\theta$ , and  $\phi_i = (\mu_i, \sigma_i)$ .

The use of  $q(\theta; \phi)$  in place of  $P(\theta|s_t, a_t)$  changes the definition of surprise from Eq. 8 to one we call variational assorted surprise for exploration (VASE):

$$U_{\text{VASE}}(s_{t+1}) = \mathbb{E}_{\theta \sim q(\cdot; \phi)} [-\log P(s_{t+1}|s_t, a_t, \theta)] - \delta H(q(\theta; \phi)). \quad (10)$$

Since the output of the model for sample  $\theta$  gives the prediction of the next state  $\hat{s}_{t+1} = f_m(s_t, a_t, \theta)$ , we define  $P(s_{t+1}|s_t, a_t, \theta)$  by measuring the deviation of  $\hat{s}_{t+1}$  from  $s_{t+1}$  under the assumption that states are normally distributed:

$$P(s_{t+1}|s_t, a_t, \theta) = \frac{1}{\sqrt{2\pi\sigma_c^2}} e^{-\|\hat{s}_{t+1} - s_{t+1}\|^2 / (2\sigma_c^2)}, \quad (11)$$

where  $\sigma_c$  is an arbitrarily chosen constant,  $\|\hat{s}_{t+1} - s_{t+1}\|$  is the norm of the difference vector between the prediction of the next state and the true next state. Note that this is a slightly different formulation to that in Equation 1 but approaches that formulation as  $\sigma_c$  approaches 0.

$N$  samples of  $\theta \sim q(\cdot; \phi)$  give  $N$  predictions for the next state from the BNN, which allows us to estimate the first term of Eq. 10 with the average:

$$\begin{aligned} & \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log P(s_{t+1}|s_t, a_t, \theta)] \\ & \approx \frac{1}{N} \sum_{n=1}^N \log P(s_{t+1}|s_t, a_t, \theta^{[n]}), \end{aligned}$$

where  $\theta^{[n]}$  is the  $n^{\text{th}}$  sample of  $\Theta$  drawn from  $q(\theta; \phi)$  and  $P(s_{t+1}|s_t, a_t, \theta^{[n]})$  is evaluated according to Eq. 11.

Since  $q(\theta; \phi)$  is a fully factorised Gaussian distribution, the second term of Eq. 10 is straight forward to evaluate:

$$\begin{aligned} H(q(\theta; \phi)) &= \sum_{i=1}^{|\Theta|} H(\mathcal{N}(\theta_i; \mu_i, \sigma_i^2)) \\ &= \frac{1}{2} \sum_{i=1}^{|\Theta|} (\log(2\pi e \sigma_i^2)). \end{aligned}$$

The last thing remaining is to ensure  $q(\theta; \phi)$  is as close as possible to  $P(\theta|\mathcal{D})$ . Variational inference uses Kullback-Leibler (KL) divergence for measuring how different  $q(\theta; \phi)$  is from  $P(\theta|\mathcal{D})$ :

$$\begin{aligned} D_{KL}[q(\theta; \phi) || P(\theta|\mathcal{D})] &= \int_{\theta} q(\theta; \phi) \log \frac{q(\theta; \phi)}{P(\theta|\mathcal{D})} \partial\theta \\ &= D_{KL}[q(\theta; \phi) || P(\theta)] \\ &\quad - \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log P(\mathcal{D}|\theta)] + \log P(\mathcal{D}). \end{aligned}$$

This difference is minimised by changing  $\phi$ , which is equivalent to maximising the variational lower bound [27]:

$$L[q(\theta; \phi), \mathcal{D}] = \mathbb{E}_{\theta \sim q(\cdot; \phi)} [\log P(\mathcal{D}|\theta)] - D_{KL}[q(\theta; \phi) || P(\theta)], \quad (12)$$

which does not require evaluation of  $P(\theta|\mathcal{D})$ . In this paper, the prior distribution of  $\theta$  is taken to be

$$P(\theta) = \prod_{i=1}^{|\Theta|} \mathcal{N}(\theta_i; 0, \sigma_m^2), \quad (13)$$

where  $\sigma_m$  is set to arbitrary value, and the expectation of log likelihood of  $P(\theta|\mathcal{D})$  is evaluated as in Eq. 12.

The entire training procedure is listed in Algorithm 1.

## IV. EXPERIMENTS AND RESULTS ANALYSIS

### A. Visualising exploration efficiency

For illustrative purposes, we begin the experimental evaluation of VASE by testing it on a simple 2DPlane environment ( $\mathcal{S} \subset \mathbb{R}^2, \mathcal{A} \subset \mathbb{R}^2$ ) which lends itself to a visualisation of the agent's exploration efficiency. The observation space is a square on the 2D plane ( $(x, y) \in \mathbb{R}^2$ ), centred on the origin. The action is its velocity  $(\dot{x}, \dot{y})$  that satisfies  $|\dot{x}| \leq 0.01, |\dot{y}| \leq 0.01$ . In this environment the agent starts at origin (0,0) and the only extrinsic reward can be found at location (1,1). The environment wraps around, so that there are no boundaries.

In this experiment, we train one agent and record the observation coordinate  $(x, y)$  in each step until it finds the non

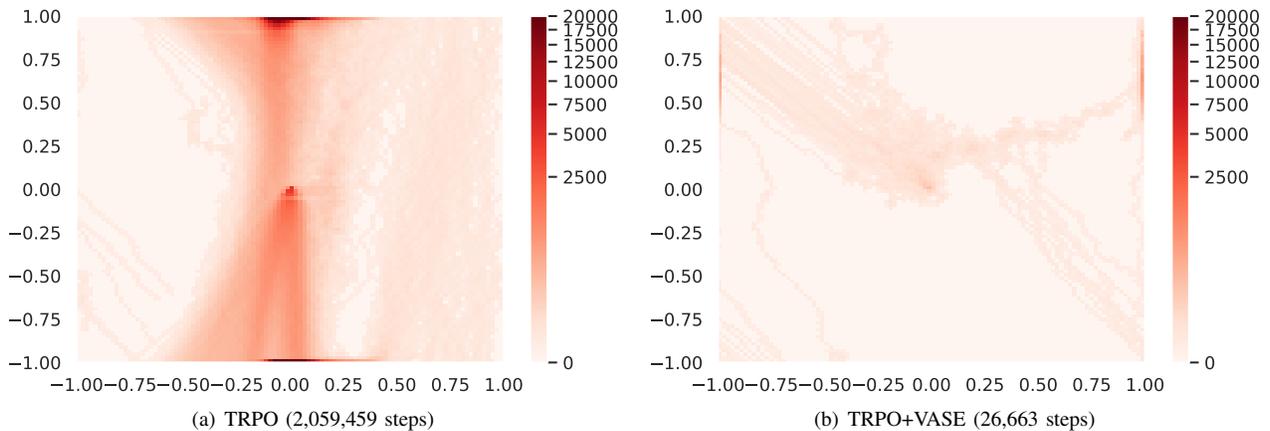


Fig. 3. Exploration efficiency as a heatmap showing the number of states visited during training in 2DPointRobot environment until chancing upon the reward state with a) no surprise, b) VASE.

---

**Algorithm 1:** Variational Assorted Surprise Exploration (VASE)

---

Initialise policy neural network  $f_\pi$  with parameters  $\psi$   
 Initialise agent’s BNN model  $f_m$ :  
     Initialise  $q(\theta; \phi)$  with parameters  $\phi$   
     Initialise prior distribution  $P(\theta)$   
 Initialise experience buffer  $R$ .  
 Reset the environment getting  $(s_0, r_e(s_0))$ .  
**for** each iteration  $n$  **do**  
     **for** each time step  $t$  **do**  
         Get action  $a_t \sim f_\pi(s_t, \psi)$   
         Sample  $\theta$   $N$  times according to  $q(\theta; \phi)$   
         Evaluate  $N$  predictions  $\hat{s}_{t+1} = f_m(s_t, a_t, \theta)$   
         Take action  $a_t$  getting  $(s_{t+1}, r_e(s_{t+1}))$   
         Compute intrinsic surprise  $U_{\text{VASE}}(s_{t+1})$   
         Construct cumulative reward  
          $r(s_{t+1}) = r_e(s_{t+1}) + \eta U_{\text{VASE}}(s_{t+1})$   
         Add new  $(s_t, a_t, s_{t+1}, r(s_{t+1}))$  to  $R$   
     **end**  
     Update  $f_m$  by maximising Eq. (12), with  $\mathcal{D}$  sampled randomly from  $R$   
     Update  $f_\pi$  using TRPO.  
**end**

---

zero extrinsic reward. Figure 3 shows the heat map of motion track for the agent trained without surprise and with VASE surprise. Darker red colour represents a higher density, which means the agent lingers more steps in this area. It is clear that random exploration strategy takes a long time (2,059,459 steps V.S. 26,663 steps) to find that first non-zero  $r_e(s_t)$ , whereas VASE does not spend time unnecessarily in random states.

**B. Continuous state/action environments**

Next we evaluate VASE on five continuous control benchmarks with very sparse reward, including three classic tasks: sparse MountainCar ( $\mathcal{S} \subset \mathbb{R}^3, \mathcal{A} \subset \mathbb{R}^1$ ), sparse Cart-PoleSwingup ( $\mathcal{S} \subset \mathbb{R}^4, \mathcal{A} \subset \mathbb{R}^1$ ), sparse Doublependulum

( $\mathcal{S} \subset \mathbb{R}^6, \mathcal{A} \subset \mathbb{R}^1$ ) and two locomotion tasks: sparse HalfCheetah ( $\mathcal{S} \subset \mathbb{R}^{20}, \mathcal{A} \subset \mathbb{R}^6$ ), sparse Ant ( $\mathcal{S} \subset \mathbb{R}^{125}, \mathcal{A} \subset \mathbb{R}^8$ ). These tasks were introduced in [7]. We also evaluate VASE on LunarLanderContinuous ( $\mathcal{S} \subset \mathbb{R}^8, \mathcal{A} \subset \mathbb{R}^2$ ) task.

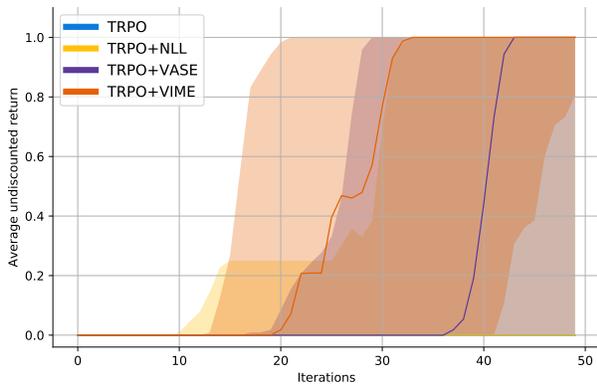
For the sparse MountainCar task, the car will climb a one-dimensional hill to reach the target. The target is located on top of a hill and on the right-hand side of the car. If the car reaches the target, the episode is done. The observation is given by the horizontal position and the horizontal velocity of the car. The agent car receives a reward of 1 only when it reaches the target.

For the sparse CartpoleSwingup task, a pole is mounted on a cart. The cart itself is limited to linear motion. Continuous cart movement is required to keep the pole upright. The system should not only be able to balance the pole, but also be able to swing it to an upright position first. The observation includes the cart position  $x$ , pole angle  $\beta$ , the cart velocity  $\dot{x}$ , and the pole velocity  $\dot{\beta}$ . The action is the horizontal force applied to the cart. The agent receives a reward of 1 only when  $\cos(\beta) > 0.9$ .

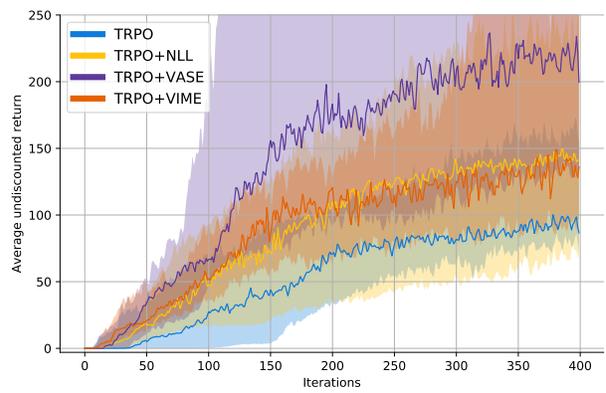
For the sparse Doublependulum task, the goal is to stabilise a two-link pendulum at the upright position. The observation includes joint angles ( $\beta_1$  and  $\beta_2$ ) and joint speeds ( $\dot{\beta}_1$  and  $\dot{\beta}_2$ ). The action is the same as in CartpoleSwingup task. The agent receives a reward of 1 only when  $\text{dist} < 0.1$ , with  $\text{dist}$  the distance between current pendulum tip position and target position.

For the sparse HalfCheetah task, the half-cheetah is a flat biped robot with nine rigid links, including two legs and one torso, and six joints. 20-dimensional observations include joint angle, joint velocity, and centroid coordinates. The agent receives a reward of 1 when  $x_{\text{body}} \geq 5$ .

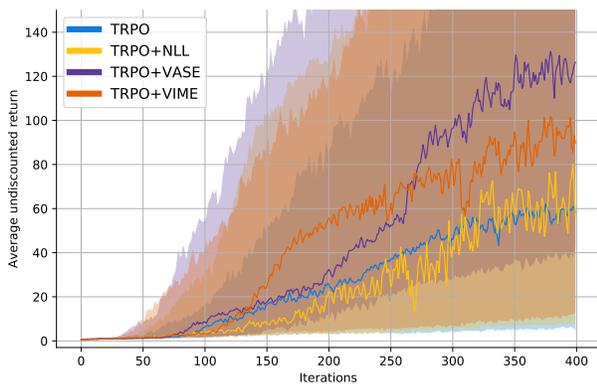
For the sparse Ant task, the ant has 13 rigid links, including four legs and a torso, along with 8 actuated joints. The 125-dim observation includes joint angles, joint velocities, coordinates of the centre of mass, a (usually sparse) vector of contact forces, as well as the rotation matrix for the body. The ant receives a reward of 1 when  $x_{\text{body}} \geq 3$ .



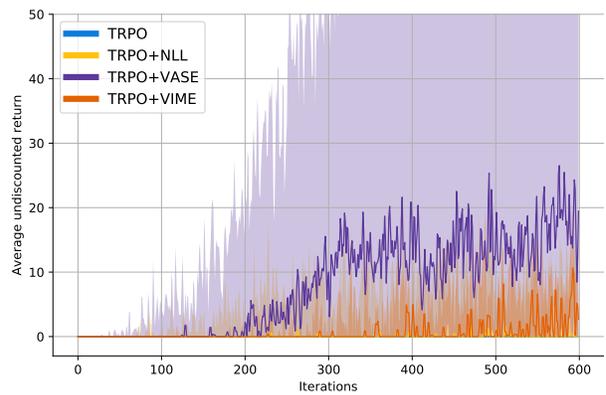
(a) MountainCar



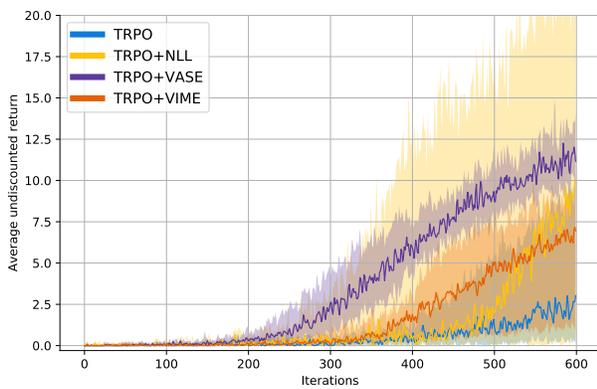
(b) CartpoleSwingup



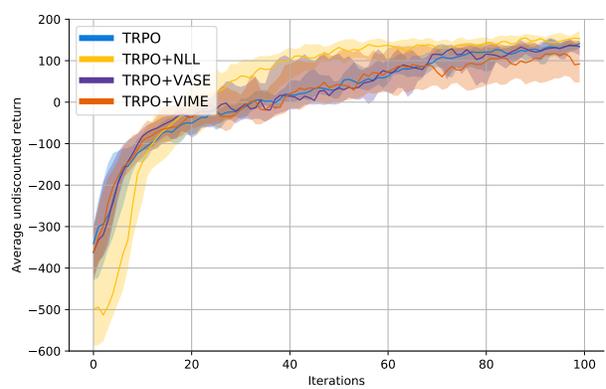
(c) DoublePendulum



(d) HalfCheetah



(e) Ant



(f) Lunar-lander

Fig. 4. Median performance for the continuous control tasks over 20 runs with a fixed set of seeds, with interquartile ranges shown in shaded areas. VIME, NLL and VASE use Bayesian surprise, surprisal, our surprise respectively.

For Lunar-lander task, the agent tries to learn to fly and then land on its landing pad. The episode is done if the lander crashes or comes to rest. The agent should get rewards of 200 when it solves this task.

All the environments except MountainCar and 2DPointRobot tasks (they are two simple tasks, we do not need to normalise them) are normalised before the algorithm starts. Here normalise the task means normalise its observations, for each observation  $o$ :

$$o = \frac{(o - \mu_o)}{\sigma_o},$$

where  $\mu_o$  and  $\sigma_o$  are the mean and standard deviation of observations. All observations and actions in these environments are continuous values. To compare with [8] and [7], we also use Trust Region Policy Optimization (TRPO) [28] method as our base reinforcement learning algorithm throughout our experiments, and we use the rllab [29] implementations of TRPO.

The number of samples drawn to compute our surprise is  $N = 10$ . The prior distribution  $P(\theta)$  is given by a Gaussian distribution from Eq. 13 with  $\sigma_m = 0.5$ .  $\sigma_c$  in Eq. 11 is set as 5. For the classic tasks sparse MountainCar, sparse CartPoleSwingup, sparse DoublePendulum and sparse LunarLanderContinuous, the  $f_m$  has one hidden layer of 32 units. All hidden layers have rectified linear unit (ReLU) non-linearities. The replay pool  $R$  has a fixed size of 100,000 samples, with a minimum size of 500 samples. For the locomotion tasks sparse HalfCheetah and sparse Ant, the  $f_m$  has two hidden layers of 64 units each. All hidden layers have tanh non-linearities. The replay pool  $R$  has a fixed size of 5,000,000 samples. The Adam learning rate of  $f_m$  is set to 0.001. All output layers are set to linear. The batch size for the policy optimisation is set to 5,000. For  $f_\pi$  the classic tasks use a neural network with one layer of 32 tanh units, while the locomotion task uses a two-layer neural network of 64 and 32 tanh units. For baseline, the classic tasks use a neural network with one layer of 32 ReLU units, while the locomotion task uses a linear function. The maximum length of trajectory LunarLanderContinuous 1000, for all the other tasks, it is 500.

Figure 4 (a)-(e) shows the median performance of three classic control tasks and two locomotion tasks. All these tasks are with sparse rewards. Figure 4 (f) shows the median performance of LunarLanderContinuous task. The agent can easily obtain rewards from this task. The performance is measured through the average return  $\mathbb{E}_\tau[\sum_{t=0}^T r_e(s_t)]$ , not including the intrinsic rewards. The median performance curves with shaded interquartile ranges areas. Figure 6 shows the speed comparison on MountainCar task.

As can be seen from Figure 4 (a)-(e), VIME performs best for sparse MountainCar task. For the sparse DoublePendulum, VIME performs well initially, but is later surpassed by VASE. VASE shows good results in sparse CartpoleSwingup, sparse HalfCheetah and sparse Ant tasks. We can also see that VASE always performs better than NLL (surprisal) in all tasks.

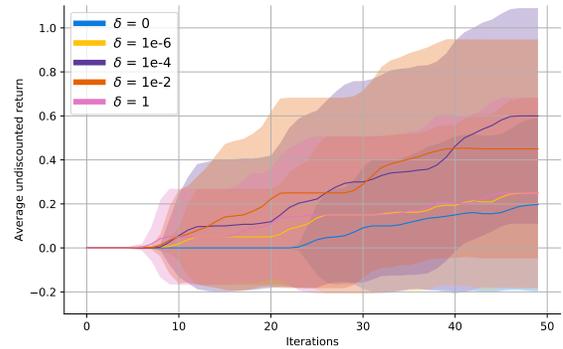


Fig. 5. Mean performance for MountainCar task with different  $\delta$ .

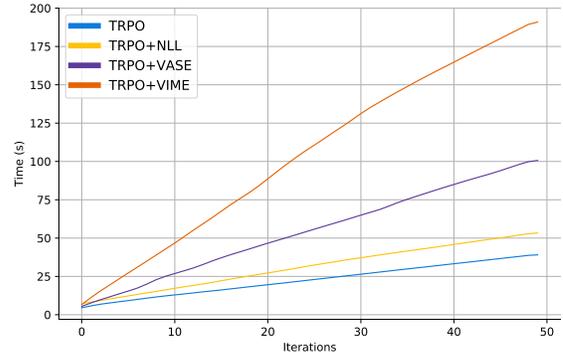


Fig. 6. Running time on MountainCar environment.

Figure 4 (f) shows that in LunarLanderContinuous task that has enough reward for the agent, all surprise-driven methods behave almost the same to the no-surprise method.

Figure 6 shows us the speed test results. For VIME, it needs to calculate a gradient through its BNN at each time step to compute the Bayesian surprise reward. This is really time consuming. However, for our VASE algorithm, it does not need to compute this gradient. Figure 6 shows that VASE runs much faster than VIME.

Finally, we also test how different values of trade-off  $\delta$  that we used in Eq. (10) affects the performs of surprise  $U_{VASE}$  on sparse MountainCar environment. We know that the value of  $H(q(M; \phi))$  depends not only on the distribution of each parameter  $M_i$ , but also on the number of parameters  $|M|$ . Meanwhile, in the beginning stages of training, the entropy of each parameter  $M_i$  is relatively large, therefore, we should take a relatively small value of  $\delta$ . Figure 5 shows VASE performance based on  $\delta$  chosen from  $\{0, 1e-6, 1e-4, 1e-2, 1\}$ . The performance is not good when  $\delta$  is too big or too small. It shows that the best interval to search  $\delta$  is  $[1e-4, 1e-2]$ .

## V. CONCLUSIONS

In this work, we chose a new form of surprise as the agent's intrinsic motivation and applied it to the RL settings by our VASE algorithm. VASE tries to approximate this surprise in a tractable way and train the agent to maximise its reward function. The agent is driven by this intrinsic reward,

which can effectively explore the environment and find sparse extrinsic rewards given by the environment. Empirical results show that VASE performs well across various continuous control tasks with sparse rewards. We believe that VASE can be easily extended to deep reinforcement learning methods or learn directly from pixel features. We leave that to future work to explore.

#### REFERENCES

- [1] Richard S. Sutton and Andrew G. Barto. *Introduction to Reinforcement Learning*. MIT press, 1998.
- [2] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharmarajan Subburaman, Daan Wierstra, Shane Legg, and Demis Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [3] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Ian Osband, Alex Graves, Vlad Mnih, Remi Munos, Demis Hassabis, Olivier Pietquin, et al. Noisy networks for exploration. *arXiv preprint arXiv:1706.10295*, 2017.
- [4] 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.
- [5] Deepak Pathak, Pulkit Agrawal, Alexei A. Efros, and Trevor Darrell. Curiosity-driven exploration by self-supervised prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 16–17, 2017.
- [6] Yuri Burda, Harri Edwards, Deepak Pathak, Amos Storkey, Trevor Darrell, and Alexei A Efros. Large-scale study of curiosity-driven learning. *arXiv preprint arXiv:1808.04355*, 2018.
- [7] Rein Houthoofd, Xi Chen, Yan Duan, John Schulman, Filip De Turck, and Pieter Abbeel. Vime: Variational information maximizing exploration. In *Advances in Neural Information Processing Systems 29*, pages 1109–1117. Curran Associates, Inc., 2016.
- [8] Joshua Achiam and Shankar Sastry. Surprise-based intrinsic motivation for deep reinforcement learning. *arXiv preprint arXiv:1703.01732*, 2017.
- [9] Alex Graves. Practical variational inference for neural networks. In *Advances in Neural Information Processing Systems 24*, pages 2348–2356. Curran Associates, Inc., 2011.
- [10] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- [11] Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimizing the description length of the weights. In *Proceedings of the Sixth ACM Conference on Computational Learning Theory*, pages 5–13. ACM Press, 1993.
- [12] Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos. Unifying count-based exploration and intrinsic motivation. In *Advances in Neural Information Processing Systems 29*, pages 1471–1479, 2016.
- [13] Georg Ostrovski, Marc G. Bellemare, Aäron van den Oord, and Rémi Munos. Count-based exploration with neural density models. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, pages 2721–2730. JMLR.org, 2017.
- [14] Manuel Lopes, Tobias Lang, Marc Toussaint, and Pierre-Yves Oudeyer. Exploration in model-based reinforcement learning by empirically estimating learning progress. In *Advances in Neural Information Processing Systems 25*, pages 206–214. Curran Associates, Inc., 2012.
- [15] Pascal Poupart, Nikos Vlassis, Jesse Hoey, and Kevin Regan. An analytic solution to discrete bayesian reinforcement learning. In *ICML*, pages 697–704. ACM, 2006.
- [16] Pierre Baldi and Laurent Itti. Of bits and wows: a bayesian theory of surprise with applications to attention. *Neural Networks*, 23(5):649–666, 2010.
- [17] Shakir Mohamed and Danilo Jimenez Rezende. Variational information maximisation for intrinsically motivated reinforcement learning. In *Advances in Neural Information Processing Systems 28*, pages 2125–2133. Curran Associates, Inc., 2015.
- [18] Jürgen Schmidhuber. A possibility for implementing curiosity and boredom in model-building neural controllers. In *Proceedings of the First International Conference on Simulation of Adaptive Behavior on From Animals to Animats*, pages 222–227. MIT Press, 1991.
- [19] Nuttapon Chentanez, Andrew G. Barto, and Satinder P. Singh. Intrinsically motivated reinforcement learning. In *Advances in Neural Information Processing Systems 17*, pages 1281–1288. MIT Press, 2005.
- [20] Myron Tribus. *Thermostatistics and thermodynamics: an introduction to energy, information and states of matter, with engineering applications*. Van Nostrand, 1961.
- [21] Jan Storck, Sepp Hochreiter, and Jürgen Schmidhuber. Reinforcement driven information acquisition in non-deterministic environments. In *Proceedings of the International Conference on Artificial Neural Networks*, volume 2, pages 159–164, 1995.
- [22] Laurent Itti and Pierre Baldi. A principled approach to detecting surprising events in video. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 631–637, 2005.
- [23] Laurent Itti and Pierre F. Baldi. Bayesian surprise attracts human attention. In *Advances in Neural Information Processing Systems 18*, pages 547–554, 2006.
- [24] Mohammadjavad Faraji. *Learning with Surprise: Theory*

*and Applications*. PhD thesis, Ecole Polytechnique Fédérale de Lausanne, 2016.

- [25] Mohammadjavad Faraji, Kerstin Preuschoff, and Wulfam Gerstner. Balancing new against old information: The role of surprise in learning. *arXiv preprint arXiv:1606.05642*, 2016.
- [26] Bradley C Stadie, Sergey Levine, and Pieter Abbeel. Incentivizing exploration in reinforcement learning with deep predictive models. *arXiv preprint arXiv:1507.00814*, 2015.
- [27] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer-Verlag, 2006.
- [28] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *ICML*, pages 1889–1897, 2015.
- [29] Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking deep reinforcement learning for continuous control. In *ICML*, pages 1329–1338, 2016.