

---

# Risk-aware Bayesian Reinforcement Learning for Cautious Exploration

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper addresses the problem of maintaining safety during training in Rein-  
2 forcement Learning (RL), such that the safety constraint violations are bounded  
3 at any point during learning. Whilst enforcing safety during training might limit  
4 the agent’s exploration, we propose a new architecture that handles the trade-off  
5 between efficient progress in exploration and safety maintenance. As the agent’s  
6 exploration progresses, we update Dirichlet-Categorical models of the transition  
7 probabilities of the Markov decision process that describes the agent’s behaviour  
8 within the environment by means of Bayesian inference. We then propose a way to  
9 approximate moments of the agent’s belief about the risk associated to the agent’s  
10 behaviour originating from local action selection. We demonstrate that this ap-  
11 proach can be easily coupled with RL, we provide rigorous theoretical guarantees,  
12 and we present experimental results to showcase the performance of the overall  
13 architecture.

## 14 1 Introduction

15 Traditionally, RL is principally concerned with the policy that the agent generates by the end of the  
16 learning process. In other words, the agent’s policy *during* learning is overlooked to the benefit of  
17 learning how to behave optimally. Accordingly, many standard RL methods rely on the assumption  
18 that the agent selects each available action at every state infinitely often during exploration [30, 28].  
19 A related technical assumption that is often made is that the MDP is *ergodic*, meaning that every  
20 state is reachable from every other state under proper action selection [25]. These assumptions may  
21 sometimes be reasonable, e.g., in virtual environments where restarting is always an option. However,  
22 in safety-critical systems, these assumptions might be unreasonable, as we may explicitly require  
23 the agent to never visit certain unsafe states. Indeed, in a variety of RL applications the safety of the  
24 agent is particularly important, e.g. expensive autonomous platforms or robots that work in proximity  
25 of humans. Thus, researchers are recently paying increasing attention not only to maximising a  
26 long-term task-driven reward, but also to enforcing avoidance of unsafe training.

27 **Related Work** The general problem of *Safe* RL has been an active area of research in which  
28 numerous approaches and definitions of safety have been proposed [3, 12, 26]. In [25], safety is  
29 defined in terms of ergodicity, with the goal of safety being that an agent is always able to return to  
30 its current state after moving away from it. In [8], safety is pursued by minimising a cost associated  
31 with worst-case scenarios, when cost is associated with a lack of safety. Similarly, [24] defines the  
32 safety constraint in terms of the expected sum of a vector of measurements to be in a target set. Other  
33 approaches [21, 16, 17, 18, 4, 19] define safety by the satisfaction of temporal logical formulae of the  
34 learnt policy, but do not provide safety *while* training such policy. Many existing approaches have  
35 been concerned with providing guarantees on the safety of the learned policy sometimes under the  
36 assumption that a backup policy is available [10, 27, 13, 22, 9, 23]. These methods are applicable to

37 systems if they can be trained on accurate simulations, but for many other real-world systems we  
 38 instead require safety *during* training.

39 There has also been much research done into the development of approaches to maintaining safety  
 40 during training. For instance, [2, 20, 15] leverage the concept of a *shield* that stops the agent from  
 41 choosing any unsafe actions. The shield assumes the agent has to observe the entire MDP (and  
 42 opponents) to construct a safety (game) model, which will be unavailable for many partially-known  
 43 MDP tasks. The approach in [11] assumes a predefined safe baseline policy that is most likely  
 44 sub-optimal, and attempts to slowly improve it with a slightly noisy action-selection policy, while  
 45 defaulting to the baseline policy whenever a measure of safety is exceeded. However, this measure  
 46 of safety assumes that nearby states have similar safety levels, which may not always be the case.  
 47 Another common approach is to use expert demonstrations to attempt to learn how to behave safely [1],  
 48 or even to include an option to default to an expert when the risk is too high [32]. Obviously, such  
 49 approaches rely heavily on the presence and help of an expert, which cannot always be counted upon.  
 50 Other approaches on this problem [35, 7, 33] are either computationally expensive or require explicit,  
 51 strong assumptions about the model of agent-environment interactions. Crucially, maintaining safety  
 52 in RL by efficiently leveraging available data is an open problem [31].

53 **Contributions** We tackle the problem of synthesising a policy via RL that optimises a discounted  
 54 reward, while not violating a safety requirement *during* learning. This paper puts forward a *cautious*  
 55 *RL scheme* that assumes the agent maintains a Dirichlet-Categorical model of the MDP. We incorporate  
 56 higher-order information from the Dirichlet distributions, in particular we compute approximations  
 57 of the (co)variances of the risk terms. This allows the agent to reason about the contribution of  
 58 epistemic uncertainty to the risk level, and therefore to make better informed decisions about how  
 59 to stay safe during learning. We show convergence results for these approximations, and propose a  
 60 novel method to derive an approximate bound on the confidence that the risk is below a certain level.  
 61 The new method adds a functionality to the agent that prevents it from taking critically risky actions,  
 62 and instead leads the agent to take safer actions whenever possible, but otherwise leaves the agent  
 63 to explore as normal. The proposed method is versatile given that it can be added on to general RL  
 64 training schemes, in order to maintain safety during learning.

## 65 2 Background

### 66 2.1 Problem Setup

67 **Definition 2.1** A finite MDP with rewards [30] is a tuple  $M = \langle Q, A, q_0, P, Re \rangle$  where  $Q =$   
 68  $\{q^1, q^2, q^3, \dots, q^N\}$  is a finite set of states,  $A$  is a finite set of actions, without loss of generality  $q_0$   
 69 is the initial state,  $P(q'|q, a)$  is the probability of transitioning from state  $q$  to state  $q'$  after taking  
 70 action  $a$ , and  $Re(q, a)$  is a real-valued random variable which represents the reward obtained after  
 71 taking action  $a$  in state  $q$ . A realisation of this random variable (namely a sample, obtained for  
 72 instance during exploration) will be denoted by  $re(q, a)$ .

73 An agent is placed at  $q_0 \in Q$  at time step  $t = 0$ . At every time step  $t \in \mathbb{N}_0$ , the agent selects an action  
 74  $a_t \in A$ , and the environment responds by moving the agent to some new state  $q_{t+1}$  according to the  
 75 transition probability distribution, i.e.,  $q_{t+1} \sim P(\cdot|q_t, a_t)$ . The environment also assigns the agent a  
 76 reward  $re(q_t, a_t)$ . The objective of the agent is to learn how to maximise the long term reward. In the  
 77 following we explain these notions more formally.

**Definition 2.2** A policy  $\pi$  assigns a distribution over  $A$  at each state:  $\pi(a|q)$  is the probability of  
 selecting action  $a$  in state  $q$ . Given a policy  $\pi$ , we can then define a state-value function

$$v_\pi(q) = \mathbb{E}^\pi \left[ \sum_{t=0}^{\infty} \gamma^t re(q_t, a_t) \mid q_0 = q \right],$$

78 where  $\mathbb{E}^\pi[\cdot]$  denotes the expected value given that actions are selected according to  $\pi$ , and  $0 < \gamma \leq 1$   
 79 is a discount factor.

80 Specifically, this means that the sequence  $q_0, a_0, q_1, a_1, \dots$  is such that  $a_n \sim \pi(\cdot|q_n)$  and  $q_{n+1} \sim$   
 81  $P(\cdot|q_n, a_n)$ . The discount factor  $\gamma$  is a pre-determined hyper-parameter that causes immediate  
 82 rewards to be worth more than rewards in the future, as well as ensuring that this sum is well-defined,  
 83 provided the standard assumption of bounded rewards. The agent's goal is to learn an optimal policy,  
 84 namely one that maximises the expected discounted return. This is actually equivalent to finding a  
 85 policy that maximises the state-value function  $v_\pi(q)$  at every state [30].

86 **Definition 2.3** A policy  $\pi$  is optimal if, at every state  $q$ ,  $v_\pi(q) = v_*(q) = \max_{\pi'} v_{\pi'}(q)$ .

87 **Definition 2.4** Given a policy  $\pi$ , we can define a state-action-value function  $v_\pi(q, a) =$   
 88  $\mathbb{E}^\pi [\sum_{t=0}^{\infty} \gamma^t r e(q_t, a_t) | q_0 = q, a_0 = a]$ , similarly to the state-value function. This allows us to  
 89 reinterpret the state-value function as  $v_\pi(q) = \sum_a v_\pi(q, a) \pi(a|q)$ , and thus we can see that an  
 90 optimal deterministic policy  $\pi$  must assign zero probability to any action  $a$  that doesn't maximise the  
 91 state-action value function.

## 92 2.2 Dirichlet-Categorical Model of the MDP

93 We consider a model for an MDP with unknown transition probabilities [14]. The transition probabili-  
 94 ties for a given state-action pair are assumed to be described by a categorical distribution over the next  
 95 state. We maintain a Dirichlet distribution over the possible values of those transition probabilities:  
 96 since the Dirichlet distribution is conjugate to the categorical distribution, we can employ Bayesian  
 97 inference to update the Dirichlet distribution, as new observations are made while the agent explores  
 98 the environment.

99 Formally, for each state-action pair  $(q^i, a)$ , we have a Dirichlet distribution  $p_a^{i1}, p_a^{i2}, \dots, p_a^{iN} \sim$   
 100  $Dir(\alpha_a^{i1}, \alpha_a^{i2}, \dots, \alpha_a^{iN})$ . The random variable  $p_a^{ij}$  represents the agent's belief about the transition  
 101 probability  $P(q^j | q^i, a)$ . At the start of learning, the agent will be assigned a prior Dirichlet distribution  
 102 for each state-action pair, according to its initial belief about the transition probabilities. At every  
 103 time step, as the agent moves from some state  $q^i$  to some state  $q^k$  by taking action  $a$ , it will generate  
 104 an event  $q^i \xrightarrow{a} q^k$ , which constitutes new data for the Bayesian inference. From Bayes' rule:

$$\begin{aligned} Pr(\mathbf{p}_a^i = \mathbf{x}_a^i | q^i \xrightarrow{a} q^k) &\propto Pr(q^i \xrightarrow{a} q^k | \mathbf{p}_a^i = \mathbf{x}_a^i) Pr(\mathbf{p}_a^i = \mathbf{x}_a^i) \\ &= x_a^{ik} \prod_j (x_a^{ij})^{\alpha_a^{ij} - 1} = \left[ \prod_{j \neq k} (x_a^{ij})^{\alpha_a^{ij} - 1} \right] (x_a^{ik})^{(\alpha_a^{ik} + 1) - 1}, \end{aligned}$$

105 which immediately yields

$$Pr(\mathbf{p}_a^i = \mathbf{x}_a^i | q^i \xrightarrow{a} q^k) = Dir(\alpha_a^{i1}, \alpha_a^{i2}, \dots, \alpha_a^{ik} + 1, \dots, \alpha_a^{iN}).$$

106 Thus, the posterior distribution is also a Dirichlet distribution. This update is repeated at each time step:  
 107 the relevant information to the agent's posterior belief about the transition probabilities is the starting  
 108 prior  $Dir(\alpha_a^{i1}, \alpha_a^{i2}, \dots, \alpha_a^{iN})$  and the "transition counts"  $c_a^{ij}$ , keeping track of the number of times that  
 109  $q^i \xrightarrow{a} q^j$  has occurred. The agent's posterior is then  $(p_a^{i1}, p_a^{i2}, \dots, p_a^{iN}) \sim Dir(\alpha_a^{i1}, \alpha_a^{i2}, \dots, \alpha_a^{iN})$ : from  
 110 this distribution, we can distill the expected value  $\bar{p}_a^{ij}$  of each random variable  $p_a^{ij}$ , as well as the  
 111 covariance of any two  $p_a^{ij}$  and  $p_a^{ik}$  (therefore also the variance of a single  $p_a^{ij}$ ):

$$\bar{p}_a^{ij} = \mathbb{E}[p_a^{ij}] = \frac{\alpha_a^{ij}}{\alpha_a^{i0}}, \quad Cov[p_a^{ij}, p_a^{ik}] = \frac{\alpha_a^{ij}(\delta^{jk} \alpha_a^{i0} - \alpha_a^{ik})}{(\alpha_a^{i0})^2 (\alpha_a^{i0} + 1)},$$

112 where  $\alpha_a^{i0} = \sum_{k=1}^N \alpha_a^{ik}$ , and  $\delta^{jk}$  is the Kronecker delta.

## 113 3 Risk-aware Bayesian RL for Cautious Exploration

114 In this section we propose a new approach to Safe RL, which will specifically address the problem of  
 115 how to learn an optimal policy in an MDP with rewards, while avoiding certain states classified as  
 116 unsafe during training. The agent is assumed to know which states of the MDP are safe and which  
 117 are unsafe, but instead of assuming that the agent has this information globally, namely for all states  
 118 of the MDP, we find it more reasonable that the agent observes states within an area around itself.  
 119 This closely resembles real-world situations, where systems may have sensors that allow them to  
 120 detect close-by danger areas, but not necessarily know about danger zones that are far away from  
 121 them. In particular, we assume that there is an observation "boundary"  $O$ , such that the agent can  
 122 observe all states that are reachable from the current state within  $O$  steps and distinguish which of  
 123 those states are safe or unsafe. The rest of this section is structured as follows:

- 124 1. In Section 3.1, we define the risk  $r_c^m(a)$  over  $m$  steps of taking an action  $a$  at the current  
 125 state, denoted as  $q^c$ . We then introduce a random variable  $R_c^m(a)$  representing the agent's  
 126 belief about the risk;
- 127 2. In Section 3.2, we leverage a method from [6] to approximate the expected value and variance  
 128 of the random variable  $R_c^m(a)$ . We provide convergence results on the approximations of  
 129 the expectation and variance of  $R_c^m(a)$ ;

- 130 3. In Section 3.3, we show how the Cantelli Inequality [5] allows us to estimate a confidence  
 131 bound on the risk  $r_c^m(a)$ ;  
 132 4. In Section 3.4, we prescribe a methodology for incorporating the expectation and variance  
 133 of risk into the action selection during the training of an RL agent.

### 134 3.1 Definition and Characterisation of the Risk

135 Given the observation boundary  $O$ , we reason about the risk incurred over the next  $m$  steps after  
 136 taking a particular action  $a$  in the current state  $q^c$ , for any  $m \leq O$ . However, note that there is a  
 137 dependence between the agent’s estimate of such a risk and the use of that estimate to inform its  
 138 action selection policy. In order to solve this dilemma we sever the dependency between the risk that  
 139 we calculate and the actions selected generating that risk by fixing a policy over the  $m$ -step horizon,  
 140 and calculating the risk given that policy. Similar to temporal-difference learning schemes, this is  
 141 done by assuming best-case action selection, namely, the  $m$ -step risk  $r_c^m(a)$  at state  $q^c$  after taking  
 142 action  $a$  is defined assuming that after selecting action  $a$ , the agent will select subsequent actions  
 143 to minimize the expected risk going forward. Assuming that the agent is at state  $q^c$ , we define the  
 144 agent’s approximation of the  $m$ -step risk  $\bar{R}_c^m(a)$  by back-propagating the risk given the “expected  
 145 safest policy” over  $m$  steps, as follows:

$$\bar{R}_k^0 = \mathbb{1}(q^k \text{ is observed and unsafe}); \quad (1)$$

$$\bar{R}_k^{n+1}(a) = \begin{cases} 1 & \text{if } q^k \text{ is observed and unsafe} \\ \sum_{j=1}^N \bar{p}_a^{kj} \bar{R}_j^n & \text{otherwise;} \end{cases} \quad (2)$$

$$\bar{R}_k^{n+1} = \begin{cases} 1 & \text{if } q^k \text{ is observed and unsafe} \\ \min_{a \in A} \bar{R}_k^{n+1}(a) & \text{otherwise.} \end{cases} \quad (3)$$

146 We terminate this iterative process at  $n + 1 = m$  and once we have calculated  $\bar{R}_c^m(a)$  ( $c = k$ ) for all  
 147 actions  $a$ . Note that, despite the use of progressing indices  $n$ , this is an iterative back-propagation  
 148 that leverages the expected values of agent’s belief about the transition probabilities, i.e.,  $\bar{p}_a^{kj}$ . Thus,  
 149  $\bar{R}_c^m(a)$  is the agent’s approximation of the expectation of the probability of entering an unsafe state  
 150 within  $m$  steps by selecting action  $a$  at state  $q^c$ , and thereafter by selecting actions that it currently  
 151 believes will minimize the probability of entering unsafe states over the given time horizon.

152 The term  $\bar{p}_a^{kj} = \mathbb{E}[p_a^{kj}]$  is used as a point estimate of the true transition probability  $t_a^{kj} = P(q^j | q^k, a)$ .  
 153 The value of  $\bar{R}_c^m$  only relies on states which the agent believes are reachable from  $q^c$  within  $m$  steps.  
 154 In particular so long as the horizon  $m$  is less than the observation boundary  $O$ , the agent is able to  
 155 observe all states which are relevant to the calculation of  $\bar{R}_c^m(a)$ , so specifically,  $\mathbb{1}(q^j \text{ is unsafe}) =$   
 156  $\mathbb{1}(q^j \text{ is observed and unsafe})$  for all relevant states  $q^j$  (see Appendix G for more details).

### 157 3.2 Approximation of Expected Value and Covariance of the Risk

158 Let  $\mathbf{x}$  denote the vector of variables  $x_a^{ij}$  where  $i, j$  range from 1 to  $N$  and  $a$  ranges over  $A$ , i.e.,  
 159  $\mathbf{x} = ((x_a^{ij})_{i,j=1,\dots,N} \text{ and } a \in A)$ . We assume that these indices are ordered lexicographically by  $(i, a, j)$ .  
 160 This is because  $i$  and  $a$  will be used to signify a state-action pair  $(q^i, a)$ , and  $j$  will be used to signify  
 161 a potential next state  $q^j$ . Introduce a set of functions (we shall see they take the shape of polynomials)  
 162  $g_k^n[\mathbf{x}]$  defined, for each state  $q^k$ , as follows:

$$g_k^0[\mathbf{x}] := \mathbb{1}(q^k \text{ is observed and unsafe});$$

$$g_k^{n+1}(a)[\mathbf{x}] := \begin{cases} 1 & \text{if } q^k \text{ is observed and unsafe} \\ \sum_{j=1}^N x_a^{kj} g_j^n[\mathbf{x}] & \text{otherwise;} \end{cases}$$

$$g_k^{n+1}[\mathbf{x}] := \begin{cases} 1 & \text{if } q^k \text{ is observed and unsafe} \\ g_k^{n+1}(\arg \min_a \bar{R}_k^{n+1}(a))[\mathbf{x}] & \text{otherwise.} \end{cases}$$

163 Then we can write the risk (of selecting action  $a$  in state  $q^c$ , over  $m$  steps) defined above as  
 164  $r_c^m(a) = g_c^m(a)[\mathbf{t}]$ , where  $\mathbf{t} = ((t_a^{ij})_{i,j=1,\dots,N} \text{ and } \forall a \in A)$  is a vector of all “true” transition prob-  
 165 abilities  $t_a^{ij} := P(q^j | q^i, a)$ . We can similarly write the agent’s approximation of the risk as  
 166  $\bar{R}_c^m(a) = g_c^m(a)[\bar{\mathbf{p}}]$ , where similarly  $\bar{\mathbf{p}} = ((\bar{p}_a^{ij})_{i,j=1,\dots,N} \text{ and } a \in A)$ . We refer to the actions spec-  
 167 ified by these argmin operators as the *agent’s expected safest action* in each state over the next  
 168  $m$  steps. Now, crucially, we can also define a new random variable  $R_c^m(a) = g_c^m(a)[\mathbf{p}]$ , where

169  $\mathbf{p} = ((p_a^{ij})_{i,j=1,\dots,N} \text{ and } \forall a \in A)$ . Since the  $p_a^{ij}$ s are random variables representing the agent's beliefs  
 170 about the true transition probabilities  $t_a^{ij}$ , we in fact have that this random variable  $R_c^m(a)$  represents  
 171 the agent's beliefs about the true risk  $r_c^m(a)$ . In the following, we show that  $\bar{R}_c^m(a)$  can be viewed as  
 172 an approximation of  $\mathbb{E}[R_c^m(a)]$ , and we provide and justify an approximation of  $Var[R_c^m(a)]$ . These  
 173 approximations can be used by the agent to reason about the true risk  $r_c^m(a)$ .

174 In order to construct approximations of the expectation and variance of  $R_c^m(a)$ , we make use of  
 175 the first-order Taylor expansion of  $g_c^m(a)[\mathbf{x}]$  around  $\mathbf{x} = \bar{\mathbf{p}}$ , following a method in [6]. The Taylor  
 176 expansion is

$$g_c^m(a)[\mathbf{x}] = g_c^m(a)[\bar{\mathbf{p}}] + \sum_{i,j=1}^N \sum_{b \in A} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} (x_b^{ij} - \bar{p}_b^{ij}) + \text{remainder}, \quad (4)$$

177 where the partial derivatives are also evaluated at  $\bar{\mathbf{p}}$ . Now we can turn equation 4 into a statistical  
 178 approximation by dropping the remainder and reasoning over the random variables  $\mathbf{p}$  for  $\mathbf{x}$ , namely:

$$g_c^m(a)[\mathbf{p}] \approx g_c^m(a)[\bar{\mathbf{p}}] + \sum_{i,j=1}^N \sum_{b \in A} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} (p_b^{ij} - \bar{p}_b^{ij}). \quad (5)$$

179 We can then take the expectation of both sides, obtaining

$$\begin{aligned} \mathbb{E}[g_c^m(a)[\mathbf{p}]] &\approx \mathbb{E}[g_c^m(a)[\bar{\mathbf{p}}]] + \mathbb{E}\left[\sum_{i,j=1}^N \sum_{b \in A} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} (p_b^{ij} - \bar{p}_b^{ij})\right] \\ &= g_c^m(a)[\bar{\mathbf{p}}] + \sum_{i,j=1}^N \sum_{b \in A} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} \mathbb{E}[(p_b^{ij} - \bar{p}_b^{ij})] = g_c^m(a)[\bar{\mathbf{p}}], \end{aligned}$$

180 where the above steps follow since the only random term in the right-hand side is  $p_b^{ij}$ , for which  
 181  $\mathbb{E}(p_b^{ij}) = \bar{p}_b^{ij}$ . Recall that  $g_c^m(a)[\mathbf{p}] = R_c^m(a)$  and  $g_c^m(a)[\bar{\mathbf{p}}] = \bar{R}_c^m(a)$ . Thus, we now have  $\bar{R}_c^m(a)$   
 182 as an approximation of the expectation of  $R_c^m(a)$ . For the approximation of the variance of the  
 183 agent's believed risk, which is again a random variable, we can write:

$$\begin{aligned} Var(g_c^m(a)[\mathbf{p}]) &\approx \mathbb{E}[(g_c^m(a)[\mathbf{p}] - g_c^m(a)[\bar{\mathbf{p}}])^2] \\ &\approx \mathbb{E}\left[\left(\sum_{i,j=1}^N \sum_{b \in A} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} (p_b^{ij} - \bar{p}_b^{ij})\right)^2\right] \quad (\text{from equation 5}) \\ &= \sum_{i,j,s,t=1}^N \sum_{b_1, b_2 \in A} \frac{\partial g_c^m(a)}{\partial x_{b_1}^{ij}} \frac{\partial g_c^m(a)}{\partial x_{b_2}^{st}} Cov(p_{b_1}^{ij}, p_{b_2}^{st}) \\ &= \sum_{i=1}^N \sum_{b \in A} \sum_{j,t=1}^N \frac{\partial g_c^m(a)}{\partial x_b^{ij}} \frac{\partial g_c^m(a)}{\partial x_b^{it}} Cov(p_b^{ij}, p_b^{it}) = \bar{V}_c^m(a), \quad (6) \end{aligned}$$

184 where  $\bar{V}_c^m(a)$  is the approximation for the variance of  $R_c^m(a)$ , i.e.,  $\approx Var(R_c^m(a))$ , and the last  
 185 line follows from the fact that the covariance between two transition probability beliefs  $p_{b_1}^{ij}$  and  $p_{b_2}^{st}$   
 186 is always 0, unless they correspond to the same starting state-action pair  $(q^i, b)$ . In other words,  
 187  $Cov(p_{b_1}^{ij}, p_{b_2}^{st}) = 0$  unless  $i = j$  and  $b_1 = b_2$ . Next, we show consistency of the estimate in the limit  
 188 (see Appendix E for the proof).

189 **Theorem 3.1** *Under Q-learning convergence assumptions [34], namely that reachable state-action*  
 190 *pairs are visited infinitely often, the estimate of the mean of the believed risk distribution  $\bar{R}_c^m(a)$*   
 191 *converges to the true risk  $r_c^m(a)$ , and it does so with the variance of the believed risk distribution*  
 192  *$Var(g_c^m(a)[\mathbf{p}])$  approaching the estimate of that variance  $\bar{V}_c^m(a)$ . Specifically,*

$$\frac{(\bar{R}_c^m(a) - r_c^m(a))}{\sqrt{\bar{V}_c^m(a)}} \rightarrow \mathcal{N}(0, 1) \text{ in distribution.}$$

193 **3.3 Estimating a Confidence on the Approximation of the Risk**

194 So far we have shown that when the agent is in the current state  $q^c$ , for each possible action  $a$ ,  
 195 approximations of the expectation and variance of its belief  $R_c^m(a)$  about the risk  $r_c^m(a)$  can be  
 196 formally obtained: we denote these two approximations by  $\bar{R}_c^m(a)$  and  $\bar{V}_c^m(a)$ , respectively. We  
 197 describe a method for combining these approximations to obtain a bound on the level of confidence  
 198 that the risk  $r_c^m(a)$  is below a certain threshold.

199 We appeal to the Cantelli Inequality, which is a one-sided Chebychev bound [5]. Having computed  
 200  $\bar{R}_c^m(a)$  and  $\bar{V}_c^m(a)$ , for a particular confidence value  $0 < C < 1$  we can define  $P_a := \bar{R}_c^m(a) +$   
 201  $\sqrt{\frac{\bar{V}_c^m(a)C}{1-C}}$ . From the Cantelli Inequality we then have

$$\Pr(R_c^m(a) \leq P_a) \geq C.$$

202 Specifically,  $P_a$  is the lowest risk level such that, according to its approximations, the agent can be at  
 203 least  $100 \times C$  % confident that the true risk is below level  $P_a$ . The agent can therefore leverage  $P_a$   
 204 when attempting to perform safe exploration (please refer to Appendix F for more details).

205 **3.4 Risk-aware Bayesian RL for Cautious Exploration (RCRL)**

206 We propose a setup for Safe RL that leverages the expectation and variance of the risk to allow an  
 207 agent to explore the environment safely, while attempting to learn an optimal policy. In order to pick  
 208 the most optimal yet safe action at each state, we propose a *double-learner* architecture, referred to  
 209 as *Risk-aware Cautious RL (RCRL)* and explained next.

210 The first learner is an optimistic agent that employs Q-learning (QL) to maximize the expected  
 211 cumulative return. The second learner is a pessimistic agent that maintains a Dirichlet-Categorical  
 212 model of the transition probabilities of the MDP. In particular, this agent is initialized with a prior  
 213  $Pri$  that encodes any information the agent might have about the transition probabilities. For each  
 214 state-action pair  $(q^i, a)$  we have a Dirichlet distribution  $p_a^{i1}, p_a^{i2}, \dots, p_a^{iN} \sim Dir(\alpha_a^{i1}, \alpha_a^{i2}, \dots, \alpha_a^{iN})$ . As  
 215 the agent explores the environment, the Dirichlet distributions are updated using Bayesian inference.

216 For each action  $a$  available in the current state  $q^c$ , the pessimistic learner computes the approximations  
 217  $\bar{R}_c^m(a)$  and  $\bar{V}_c^m(a)$  of its belief  $R_c^m(a)$  of the risk over the next  $m$  steps of taking action  $a$  in  $q^c$ . The  
 218 “risk horizon”  $m$  is a hyper-parameter that, as discussed, should be set at most as the observation  
 219 boundary  $O$ . The pessimistic learner is also initialized with two hyper-parameters  $P_{max}$  and  $C(n)$ :  
 220  $P_{max}$  represents the maximum level of risk that the agent should be prepared to take, whereas  $C(n)$  is  
 221 a decreasing function of the number of times  $n$  that the current state has been visited, which satisfies  
 222  $C(0) < 1$  and  $\lim_{n \rightarrow \infty} C(n) = 0$ . From Section 3.3, the agent can then compute, for each action  $a$ ,  
 223 the value

$$P_a = \bar{R}_c^m(a) + \sqrt{\frac{\bar{V}_c^m(a)C(n)}{1-C(n)}}, \quad (7)$$

224 which can thus define a set of safe actions: these are all the actions that the agent believes have risk  
 225 less than  $P_{max}$ , with confidence at least  $C(n)$ , namely

$$A_{safe} = \{a \in A | P_a \leq P_{max}\}.$$

226 In case there are no actions  $a$  such that  $P_a \leq P_{max}$ , the agent instead allows

$$A_{safe} = \{a \in A | \bar{R}_c^m(a) = \min_{a'} \bar{R}_c^m(a)\}. \quad (8)$$

227 Finally, the agent selects an action  $a^*$  from the set of safe actions using softmax action selection [30]  
 228 according to the Q-values of those actions, with some *temperature*  $t > 0$ :

$$Pr(a^* = a) = \frac{e^{Q(q^c, a)/t}}{\sum_{a \in A_{safe}} e^{Q(q^c, a)/t}}. \quad (9)$$

229 The pseudo-code for the full algorithm is available in Appendix B.

230 In summary, we effectively have two agents learning to accomplish two tasks. The first agent  
 231 performs Q-learning to learn an optimal policy for the reward. The second agent determines the  
 232 best approximation of the expected value and variance of each action, enabling it to prevent the  
 233 first agent from selecting actions that it cannot guarantee to be safe enough (with at least a given

Table 1: Total successes and failures. Gridworld: different priors and acceptable risks  $P_{\max}$ , averaged over 10 agents. PacMan: varying risk horizon  $m$ , single agent.

Experiment	Setup	# Successes	# Failures	Total Episodes
Gridworld	Prior 1, $P_{\max} = 0.33$	404.3	54.2	500
	Prior 1, $P_{\max} = 0.01$	506.0	417.9	1500
	Prior 2, $P_{\max} = 0.01$	384.6	0.5	500
	Prior 3, $P_{\max} = 0.01$	407.4	14.4	500
	Prior 3, $P_{\max} = 0.0033$	421.3	1.1	500
	Native Q-Learning	414.6	990.5	1500
PacMan	Risk Horizon $m = 2$	234	77	311
	Risk Horizon $m = 3$	207	68	275
	Native Q-Learning	0	1500	1500

confidence). When instead the pessimistic agent cannot guarantee that any action is safe enough, it forces the optimistic learner to go into “safety mode”, i.e., to forcibly select the actions that minimize the expected value of the risk, as per equation 8. From an empirical perspective, implementing this concept of a “safety mode” allows for continued progress, and pairs extremely well with the definition of the risk: namely, when the agent deems that a state is too risky, it will go into this “safety mode” until it is back in a state with sufficiently safe actions.

Finally, note that  $C(n)$  represents the level of confidence that the agent requires in an action being safe enough for it to consider taking that action. When the agent starts exploring and  $C(n)$  is at its highest, the agent only explores actions that it is very confident in. However, it may need to take actions that it is less confident in order to find an optimal policy. Thus, as it continues exploring,  $C(n)$  is reduced, allowing the agent to select actions upon which it is not as confident. However, in the limit, when  $C(n) \rightarrow 0$ , we have that  $P_a = \bar{R}_c^m(a)$ , which means that the agent never takes an action if its approximation of the expected value of the risk  $\bar{R}_c^m(a)$  is more than the maximum allowable risk  $P_{\max}$ .

## 4 Experiments

**Gridworld** - We first evaluated the performance of RCRL on a *Slippery Gridworld Bridge* example. The states of the MDP consist of a  $20 \times 20$ -grid, as depicted in Figure 2a (Appendix C). The agent is initialized at  $q_0$  in the bottom-left corner (green). The agent’s task is to get to the goal region without ever entering an unsafe state. In particular, upon reaching a goal state, the agent is given a reward of 1 and the learning episode is terminated; at every other state it receives a reward of 0, and upon reaching an unsafe state the learning episode terminates with reward 0. At each time step the agent might move into one of the 4 neighbouring states, or stay in its current position; thus, the agent has access to 5 actions at each state,  $A = \{right, up, left, down, stay\}$ . If the agent selects action  $a \in A$ , then it has a 96% chance of moving in direction  $a$ , and a 4% chance of “slipping”, namely moving to another random direction. If any movement would ever take the agent outside of the grid, then the agent will just remain in place. The agent is assumed to have an observation boundary  $O = 2$  steps. Note that due to the slipperiness of the grid and the narrow passage to reach the goal state, minimizing the risk is not aligned with maximizing the expected reward.

We tested RCRL with 5 different combinations of a prior  $Pri$  and a maximum acceptable risk  $P_{\max}$ . The following additional hyper-parameters of the algorithm were kept constant: the maximum number of steps per episode  $max\_steps = 400$ , the maximum number of episodes  $max\_episodes = 500$  (although this was increased to 1500 in two cases when the agent did not converge to near-optimal policy within the first 500, cf. Table 1); the learning rate  $\mu = 0.85$ ; the discount factor  $\gamma = 0.9$ ; and the risk horizon  $m = 2$  (Appendix B). Recall that a prior consists of a Dirichlet distribution  $p_a^{i1}, \dots, p_a^{iN} \sim Dir(\alpha_a^{i1}, \dots, \alpha_a^{iN})$  for every state-action pair  $(q^i, a)$ . We considered three priors:

- Prior 1 - completely uninformative: in this case we assigned a value of 1 to every  $\alpha$ . This yields a distribution that is uniform over its support.
- Prior 2 - weakly informative: we assigned a value of 12 to the  $\alpha$  corresponding to moving in the correct direction, and a value of 1 to all other  $\alpha$ ’s. This gives a distribution in between Prior 1 and Prior 3 in both degree of bias and concentration.
- Prior 3 - highly informative: we assigned a value of 96 to the  $\alpha$  corresponding to moving in the correct direction, and a value of 1 to all other  $\alpha$ ’s. This gives a distribution that is highly

276 concentrated, and for which the mean values of the transition probability random variables  
277 are the true transition probabilities of the MDP, and hence unbiased.

278 We tested the algorithm with all three priors and a maximum acceptable risk of  $P_{\max} = 0.01$  and  
279 repeating each experiment 10 times to take averages. On average, the agent with the highly informative  
280 prior (Prior 3) entered unsafe states 14.4 times, and always converged to near-optimality within  
281 about 200 steps, successfully crossing the bridge 407.4 times. For the other 78.2 episodes, the agent  
282 reached the episode limit within crossing the bridge or entering an unsafe state. The agent with Prior  
283 2 interestingly only entered unsafe states an average of 0.5 times per experiment, and converged to a  
284 near-optimal policy within about 300 episodes, successfully crossing the bridge 384.6 times. On the  
285 other hand, the agent with Prior 1 only crossed the bridge less than 30 times. We therefore increased  
286 the total number of episodes to 1500 and tried again, yet still over half the time it did not converge to  
287 a near-optimal policy (Appendix A).

288 We then tested Prior 1 with a more lenient maximum acceptable risk of  $P_{\max} = 0.33$ , and found that  
289 the agent this time managed to converge to near-optimality within around 200 episodes, entering  
290 unsafe states 54.2 times and successfully crossing the bridge 404.3 times. We also tested Prior 3 with  
291 a stricter  $P_{\max} = 0.0033$  and found out that it entered unsafe states only 1.1 times and succeeded  
292 421.3 times, converging to near-optimality within 150 episodes (Appendix A).

293 Finally, we tested native Q-learning, without any safe learning scheme. This native scheme had  
294 almost no successful crossings of the bridge in the first 500 episodes, so we ran it for 1500 episodes  
295 and found that it only converged to a near-optimal policy about half the time, on average entering  
296 unsafe states 990.5 times and successfully crossing the bridge 414.6 times.

297 Table 1 summarizes the number of successes and failures for each agent. To understand better the  
298 rate of convergence to near-optimality, Figure 1 (Appendix A) displays the number of steps taken by  
299 the agent to cross the bridge at every successful episode (it displays 400 if the agent never managed  
300 to cross the bridge) averaged over the 10 experiments. On each graph we display for comparison  
301 the theoretical least number of steps it could cross the bridge in, which is 22. Note that because the  
302 grid-world is slippery, even an optimal policy would have fluctuations above the 22-steps line.

303 **Discussion** The first result of note is how poorly Prior 1 performs with  $P_{\max} = 0.01$ . It mostly fails  
304 to converge to near-optimal behaviour even with 1500 steps as can be seen in Figure 1b (Appendix A),  
305 in fact seeming to converge slower than native Q-learning. This occurs because the maximum  
306 allowable risk is set too low for the given prior. In particular, there are two main issues with this.  
307 The first issue is a type of degenerate behaviour specific to our algorithm and to the completely  
308 uninformative prior with overly strict  $P_{\max}$ : given that the agent starts with no information on the  
309 transition probabilities, it is unable to tell which actions are safe and which are unsafe. In particular,  
310 with  $P_{\max}$  at 1%, the first time the agent arrives at any state  $q^c$  from which it can observe some  
311 unsafe state, it immediately goes into safety mode as it judges that the risk of every action is above  
312 1%. Since it has no information on which action is safest, it randomly selects an action (assuming the  
313 Q-values were initialized to 0). If that randomly-selected action does not take the agent closer to a  
314 risky state, then after updating the agent’s beliefs about the transition probabilities for that action, it  
315 will believe that action is the safest one from that state. Thus every time it encounters that state again,  
316 it will *always* select that action, never attempting any other actions. This behaviour can be seen in  
317 Figure 2b (Appendix C). The state (13, 1) has been visited significantly more often than any other  
318 state. This has occurred because the first time the agent encountered that state, it chose action *stay*,  
319 and as above, from then on always chose *stay* in state (13, 1). This would cause the agent to remain  
320 in (13, 1) until it slipped off of that state.

321 The second issue with having such a strict  $P_{\max}$  could involve any prior. In this case  $P_{\max}$  is set  
322 so low that actions that may be optimal are simply never tested, as the agent’s initial belief about  
323 those actions causes the expected risk associated with them to always be greater than  $P_{\max}$ . This  
324 should not be viewed as an undesirable consequence of the algorithm, but rather as the algorithm  
325 working as intended. With the maximum allowable risk level  $P_{\max}$  set so low, the agent judges that  
326 certain actions are riskier than acceptable and therefore does not take them. However, this does raise  
327 a more general question about the nature of safe learning in general: ensuring safety while learning  
328 necessarily means avoiding actions we believe are too dangerous, so if we want any guarantees on  
329 safety, then we must accept that the agent may be unable to explore the entire state space.

330 The second result of note is that Prior 3 performs much less safely than Prior 2 does at  $P_{\max} = 0.01$ .  
331 This seems counter intuitive at first, given that Prior 3 is more accurate and more confident than Prior  
332 2. However, the explanation is quite simple. Prior 3 (initially) causes the agent’s expected belief to  
333 correctly predict that there is only a 1% chance of moving to an unsafe state on a particular step if the



334 agent selects the action to move away from it. On the other hand, Prior 2 causes the agent’s expected  
 335 belief to predict there is a 6.25% chance of this happening. Thus, Prior 3 (correctly) evaluates the  
 336 risk of moving within 1 step of a risky state as much lower than Prior 2 does. It is likely that at some  
 337 points in the experiments, the agent with Prior 3 chose to move within 1 step of an unsafe state where  
 338 an agent starting with Prior 2 (with the same experiences) would have rejected that action as too risky.  
 339 The agent with Prior 3 would then be at risk of slipping into an unsafe state. In Figure 2c and 2d  
 340 (Appendix C), we can see exactly this happening, where Prior 3 regularly visits state (13, 8), which is  
 341 adjacent to the unsafe state (12, 8). Prior 2 instead regularly moves one more state to the right before  
 342 moving up to row 13, since (12, 9) is safe.

343 Prior 3 with  $P_{\max} = 0.0033$  shows how we can make use of a highly accurate prior to guarantee even  
 344 less risk, and in this case the agent almost never enters unsafe states, while converging faster than any  
 345 other setup to near-optimality.

346 The final result is that the rate of convergence of the native Q-learning agent is much slower on this  
 347 MDP than the other agents (excluding Prior 1 with the inappropriate  $P_{\max} = 0.01$ ). As in Figure  
 348 1 (Appendix A), Q-learning took between 300 and 1500 episodes to converge when it did, and  
 349 occasionally failed to converge, compared to 150-300 episodes for the four other agents to converge  
 350 in all 10 experiments. This was even the case for the agent with the completely uninformative prior,  
 351 with  $P_{\max} = 0.33$ . This is a key result: it shows that not only can RCRL keep the agent safe during  
 352 learning when possible, it may also direct the agent to explore more fruitful areas of the state-space.  
 353 In this case study in particular, the native Q-learning agent entered unsafe states so often initially that  
 354 it took many episodes before it was able to access the bridge and find the reward at the other side.  
 355 Conversely, since the safe agents mostly avoided “sinking” situations, they were able to explore much  
 356 more of the state space on each episode.

357 **PacMan** - We also evaluated the performance of RCRL on a *PacMan* example. Figure 3a (Ap-  
 358 pendix D) depicts the initial state of the environment, where the agent (PacMan) must get to both  
 359 yellow dots (food) without getting caught by the ghost. Note that because both the agent and the  
 360 ghost move through the maze, the PacMan MDP has about 10 times more states than the Gridworld,  
 361 and up to 5 times more possible next states at any given state. Upon picking up the second piece of  
 362 food, the agent is given a reward of 1 and the learning episode stops. Every other state incurs a reward  
 363 of 0 and if the ghost catches PacMan, the learning episode stops with reward 0. The agent has access  
 364 to four actions at each state,  $A = \{right, up, left, down\}$  and will move in the direction selected,  
 365 or if that direction moves into a wall, then it will stay still. The ghost will with 90% probability  
 366 move in the direction that takes it closest to the agent’s next location, and with 10% probability  
 367 will move in a random direction. For this setup, we assumed an observation boundary  $O = 3$  and  
 368 compared two values of the risk horizon,  $m = 2, 3$ . We therefore kept constant the other parameters  
 369 and hyper-parameters: the learning rate  $\mu = 0.85$ ; the discount factor  $\gamma = 0.9$ ; the maximum number  
 370 of steps per episode  $max\_steps = 400$ ; the maximum acceptable risk  $P_{max} = 0.33$ ; the prior, which  
 371 we set to be a completely uninformative prior as in the Gridworld example; the maximum number of  
 372 episodes, which we set as 1500 or the number of episodes before the total rate of successful episodes  
 373 exceeded 75%.

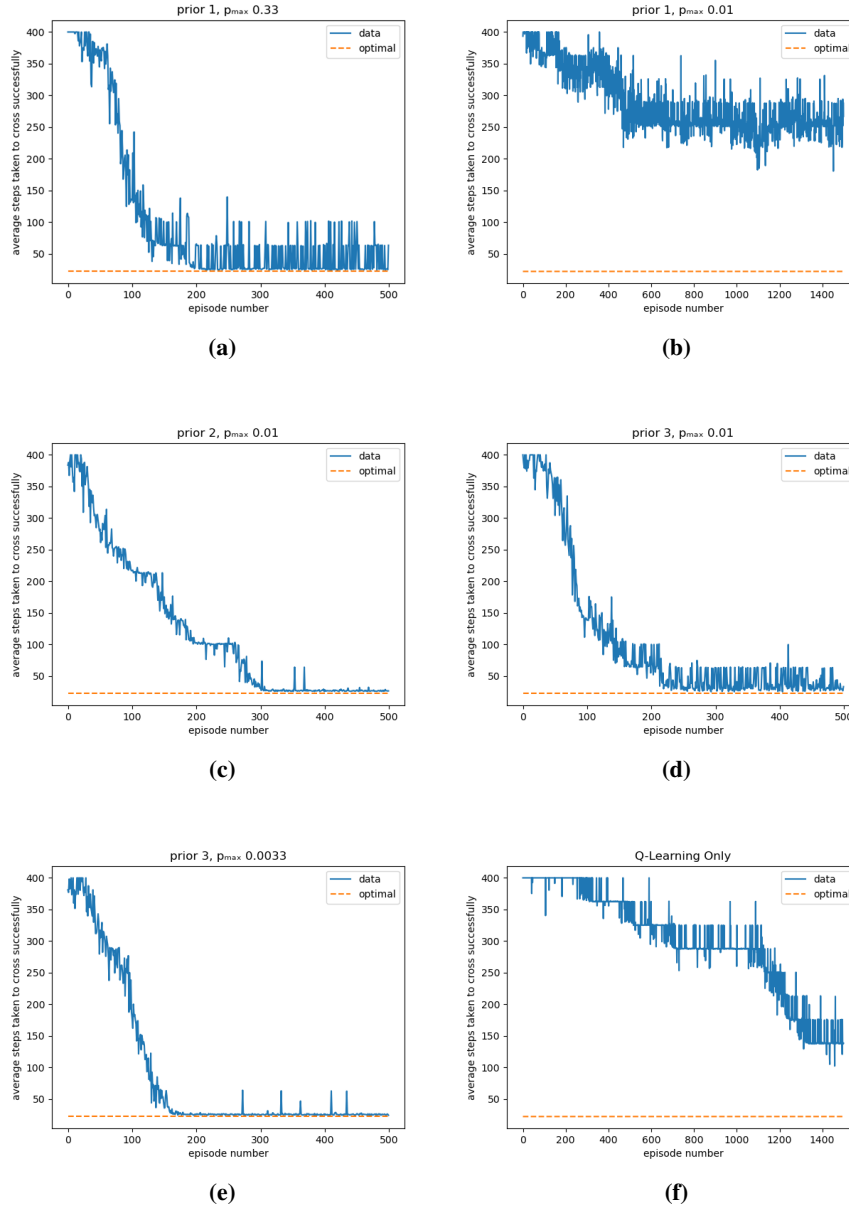
374 As in Table 1, the agent with a risk horizon of  $m = 2$  steps exceeded a success rate of 75% after 311  
 375 episodes, having failed 77 times. The agent with the larger risk horizon of  $m = 3$  only took 275 steps  
 376 to exceed that success rate, and only failed 68 times. Figures 3b-3c (Appendix D) display the number  
 377 of steps taken by the agent to win (or 400 if they lose) for each agent, as well as the running average  
 378 number of steps over the previous 50 episodes.

379 **Discussion** The improvement in performance from  $m = 2$  to 3 is likely due to the increased foresight  
 380 of the agent leading it to move away from excessively risky scenarios further in advance, potentially  
 381 avoiding entering a state from which entering a dangerous state is unavoidable. However, it may also  
 382 be simply due to the fact that increasing the risk horizon leads to an overall increase in risk estimates,  
 383 which will naturally cause more actions to be considered too risky and may reduce the number of  
 384 failures. In other words, we may have been in a situation where decreasing the maximum acceptable  
 385 risk  $P_{max}$  would have led to similar improvements, and the increase in risk horizon was behaving  
 386 functionally more like a decrease in  $P_{max}$ . Both risk-aware agents compare very favourably against  
 387 the Native Q-Learning agent, which did not succeed once in 1500 episodes.

388 **References**

- 389 [1] Pieter Abbeel, Adam Coates, and Andrew Y. Ng. Autonomous helicopter aerobatics through  
390 apprenticeship learning. *The International Journal of Robotics Research*, 29(13):1608–1639,  
391 2010.
- 392 [2] Mohammed Alshiekh, Roderick Bloem, Ruediger Ehlers, Bettina Könighofer, Scott Niekum,  
393 and Ufuk Topcu. Safe reinforcement learning via shielding, 2017.
- 394 [3] Lukas Brunke, Melissa Greeff, Adam W. Hall, Zhaocong Yuan, Siqi Zhou, Jacopo Panerati, and  
395 Angela P. Schoellig. Safe learning in robotics: From learning-based control to safe reinforcement  
396 learning. *CoRR*, abs/2108.06266, 2021.
- 397 [4] Mingyu Cai, Hosein Hasanbeig, Shaoping Xiao, Alessandro Abate, and Zhen Kan. Modular  
398 deep reinforcement learning for continuous motion planning with temporal logic. *IEEE Robotics  
399 and Automation Letters*, 6(4):7973–7980, 2021.
- 400 [5] Francesco Paolo Cantelli. Sui confini della probabilita. In *Atti del Congresso Internazionale dei  
401 Matematici: Bologna del 3 al 10 de settembre di 1928*, pages 47–60, 1929.
- 402 [6] George Casella and Roger L. Berger. *Statistical inference*. Brooks/Cole Cengage Learning,  
403 2021.
- 404 [7] Richard Cheng, Gábor Orosz, Richard M Murray, and Joel W Burdick. End-to-end safe  
405 reinforcement learning through barrier functions for safety-critical continuous control tasks. In  
406 *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 3387–3395,  
407 2019.
- 408 [8] Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-Constrained  
409 Reinforcement Learning with Percentile Risk Criteria. *Journal of Machine Learning Research  
410 18*, pages 1–51, 2018.
- 411 [9] Yinlam Chow, Ofir Nachum, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. A  
412 lyapunov-based approach to safe reinforcement learning. In *Advances in neural information  
413 processing systems*, pages 8092–8101, 2018.
- 414 [10] Stefano P. Coraluppi and Steven I. Marcus. Risk-sensitive and minimax control of discrete-time,  
415 finite-state Markov decision processes. *Automatica*, 35(2):301–309, 1999.
- 416 [11] J. Garcia and F. Fernandez. Safe exploration of state and action spaces in reinforcement learning.  
417 *Journal of Artificial Intelligence Research*, 45:515–564, Dec 2012.
- 418 [12] Javier Garcia and Fernando Fernandez. A Comprehensive Survey on Safe Reinforcement  
419 Learning. *Journal of Machine Learning Research 16*, 2015.
- 420 [13] P. Geibel and F. Wysotzki. Risk-sensitive reinforcement learning applied to control under  
421 constraints. *Journal of Artificial Intelligence Research*, 24:81–108, Jul 2005.
- 422 [14] Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, and Aviv Tamar. Bayesian Reinforce-  
423 ment Learning: A Survey. *Foundations and Trends in Machine Learning*, 8(5-6):359–483, 2015.  
424 arXiv: 1609.04436.
- 425 [15] Mirco Giacobbe, Hosein Hasanbeig, Daniel Kroening, and Hjalmar Wijk. Shielding Atari  
426 games with bounded prescience. In *Proceedings of the 20th International Conference on  
427 Autonomous Agents and MultiAgent Systems*. International Foundation for Autonomous Agents  
428 and Multiagent Systems, 2021.
- 429 [16] Hosein Hasanbeig, Alessandro Abate, and Daniel Kroening. Certified reinforcement learning  
430 with logic guidance. *arXiv preprint arXiv:1902.00778*, 2019.
- 431 [17] Hosein Hasanbeig, Alessandro Abate, and Daniel Kroening. Logically-constrained neural fitted  
432 Q-iteration. In *Proceedings of the 18th International Conference on Autonomous Agents and  
433 MultiAgent Systems*, pages 2012–2014. International Foundation for Autonomous Agents and  
434 Multiagent Systems, 2019.

- 435 [18] Hosein Hasanbeig, Daniel Kroening, and Alessandro Abate. Deep reinforcement learning  
436 with temporal logics. In *International Conference on Formal Modeling and Analysis of Timed*  
437 *Systems*, pages 1–22. Springer, 2020.
- 438 [19] Hosein Hasanbeig, Daniel Kroening, and Alessandro Abate. LCRL: Certified policy synthesis  
439 via logically-constrained reinforcement learning. In *International Conference on Quantitative*  
440 *Evaluation of Systems*, pages 217–231. Springer, 2022.
- 441 [20] Nils Jansen, Bettina Könighofer, Sebastian Junges, Alexandru C. Serban, and Roderick Bloem.  
442 Safe reinforcement learning via probabilistic shields, 2019.
- 443 [21] Xiao Li and Calin Belta. Temporal logic guided safe reinforcement learning using control  
444 barrier functions, 2019.
- 445 [22] Tommaso Mannucci, Erik-Jan van Kampen, Cornelis De Visser, and Qiping Chu. Safe explo-  
446 ration algorithms for reinforcement learning controllers. *IEEE transactions on neural networks*  
447 *and learning systems*, 29(4):1069–1081, 2017.
- 448 [23] Hongzi Mao, Malte Schwarzkopf, Hao He, and Mohammad Alizadeh. Towards safe online rein-  
449 forcement learning in computer systems. In *33rd conference on neural information processing*  
450 *systems (NeurIPS 2019)*, 2019.
- 451 [24] Sobhan Miryoosefi, Kianté Brantley, Hal Daume III, Miro Dudik, and Robert E Schapire.  
452 Reinforcement learning with convex constraints. *Advances in Neural Information Processing*  
453 *Systems*, 32, 2019.
- 454 [25] Teodor Mihai Moldovan and Pieter Abbeel. Safe Exploration in Markov Decision Processes.  
455 *arXiv:1205.4810 [cs]*, July 2012. arXiv: 1205.4810.
- 456 [26] Martin Pecka and Tomas Svoboda. Safe exploration techniques for reinforcement learning—an  
457 overview. In *International Workshop on Modelling and Simulation for Autonomous Systems*,  
458 pages 357–375. Springer, 2014.
- 459 [27] Theodore J Perkins and Andrew G Barto. Lyapunov design for safe reinforcement learning.  
460 *Journal of Machine Learning Research*, 3(Dec):803–832, 2002.
- 461 [28] Martin L. Puterman. *Markov decision processes: Discrete stochastic dynamic programming*.  
462 John Wiley & Sons, 2014.
- 463 [29] E. Slutsky. Über Stochastische Asymptoten und Grenzwerte. *Metron* 5, Nr. 3, 3-89 (1925).,  
464 1925.
- 465 [30] Richard S. Sutton, Francis Bach, and Andrew G. Barto. *Reinforcement Learning: An Introduc-*  
466 *tion*. MIT Press Ltd, 2018.
- 467 [31] Andrew J Taylor, Victor D Dorobantu, Sarah Dean, Benjamin Recht, Yisong Yue, and Aaron D  
468 Ames. Towards robust data-driven control synthesis for nonlinear systems with actuation  
469 uncertainty. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 6469–6476.  
470 IEEE, 2021.
- 471 [32] Lisa Torrey and Matthew E. Taylor. Help an agent out: Student/teacher learning in sequential  
472 decision tasks. In *Proceedings of the Adaptive and Learning Agents workshop (at AAMAS-12)*,  
473 June 2012.
- 474 [33] Matteo Turchetta, Felix Berkenkamp, and Andreas Krause. Safe exploration in finite Markov  
475 decision processes with Gaussian processes. In *Advances in Neural Information Processing*  
476 *Systems*, pages 4312–4320, 2016.
- 477 [34] Christopher John Cornish Hellaby. Watkins. *Learning from delayed rewards*. PhD thesis, 1989.
- 478 [35] Min Wen and Ufuk Topcu. Constrained cross-entropy method for safe reinforcement learning.  
479 In *Advances in Neural Information Processing Systems*, pages 7450–7460, 2018.



**Figure 1:** The number of steps it takes the agent to cross the bridge for every episode where it crosses. Averaged over 10 experiments. Results for Q-learning only and for RCRL across different priors and values of risk  $P_{\max}$ . As Q-learning converges, it approaches the lower bound on the optimal number of steps per episode.

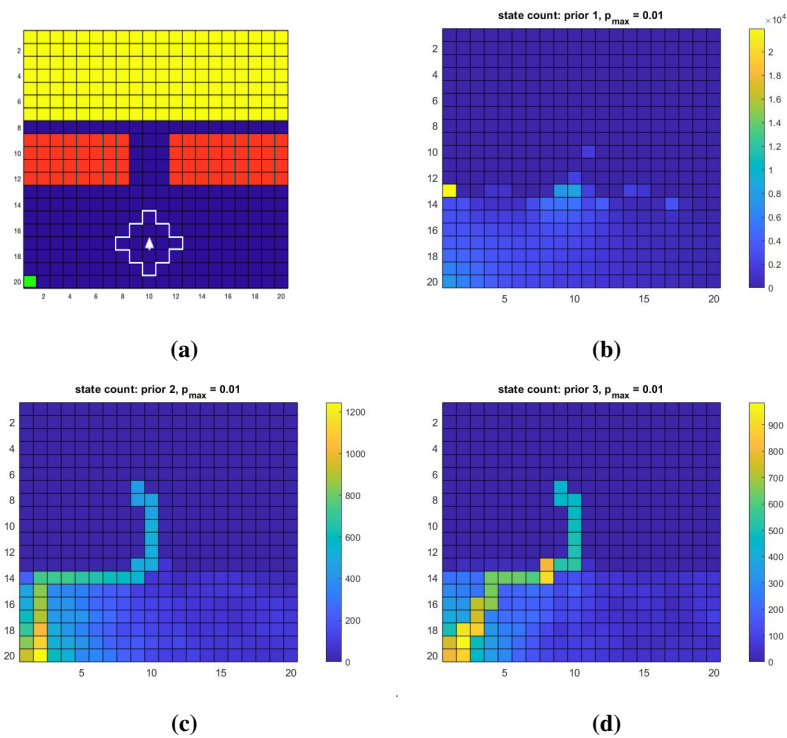
---

**Algorithm 1:** Risk-aware Cautious RL (RCRL)

---

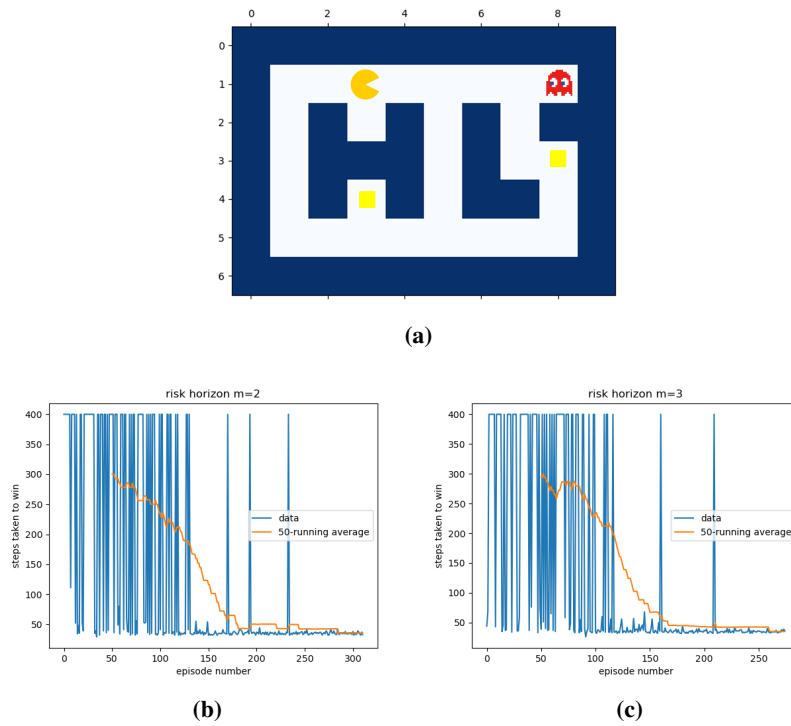
**input :**  $P, \mathcal{C}(n), P_{\max}, \max\_steps, \max\_episodes, \mu, \gamma, m$ 

- (1) initialize  $Q(q, a)$  for each state-action pair  $(q, a)$ ;
  - (2) initialize  $num\_steps = 0$ ;
  - (3) initialize  $num\_episodes = 0$ ;
  - while**  $num\_episodes < \max\_episodes$  **do**
  - (4)  $q^c \leftarrow q^0$ ;
  - (5)  $num\_episodes \leftarrow num\_episodes + 1$ ;
  - while**  $num\_steps < \max\_steps$  and  $q^c$  is not unsafe **do**
  - (6) calculate  $\bar{R}_c^m(a)$  as in (2);
  - (7) calculate  $\bar{V}_c^m(a)$  as in (6);
  - (8) calculate  $P_a$  as in (7);
  - (9)  $A_{safe} := \{a \in A | P_a \leq P_{\max}\}$ ;
  - if**  $A_{safe} = \emptyset$  **then**
  - (10)  $A_{safe} \leftarrow \{a \in A | \bar{R}_c^m(a) = \min_{a'} \bar{R}_c^m(a)\}$ ;
  - end**
  - (11) choose action  $a^*$  according to (9);
  - (12) pass action  $a^*$  to environment and receive next state  $q'$  and reward  $re(q^c, a^*)$ ;
  - (13) update belief  $p$  as in section 2;
  - (14) update  $Q(q^c, a^*) \leftarrow (1 - \mu)Q(q^c, a^*) + \mu(re(q^c, a^*) + \gamma \max_{a'} Q(q', a'))$ ;
  - (15)  $q^c \leftarrow q'$ ;
  - (16)  $num\_steps \leftarrow num\_steps + 1$ ;
  - end**
  - end**
-



**Figure 2:** (a) Slippery Gridworld setup: agent is represented by an arrow surrounded by the observation area (white line). Labels denote target (yellow), unsafe (red) and safe states (blue), and initial state ( $q_0$ , green). (b) For a single experiment, number of state-visitations for Prior 1 at  $P_{\max} = 0.01$ . (c-d) Number of state-visitations, for Priors 2 and 3 at  $P_{\max} = 0.01$ .

485 **D Appendix D. PacMan Experiment details**



**Figure 3:** (a) PacMan Setup: agent (PacMan) starts at position (1,3). Food is denoted by yellow dots, and the ghost starts in the top right corner. (b-c) Number of steps taken to win (i.e. eat both foods without being caught by the ghost) on episodes where the agent does win (or 400 if the agent is caught), for risk horizon 2 and 3. The orange line denotes the running average number of steps to win over the previous 50 episodes.

486 **E Appendix E. Convergence Results for the Approximations of the Expected**  
 487 **Value and Variance of the Risk**

488 **Theorem E.1** Under  $Q$ -learning convergence assumptions [34], namely that reachable state-action  
 489 pairs are visited infinitely often, the estimate of the mean of the believed risk distribution  $\bar{R}_c^m(a)$   
 490 converges to the true risk  $r_c^m(a)$ , and it does so with the variance of the believed risk distribution  
 491  $Var(g_c^m(a)|\mathbf{p})$  approaching the estimate of that variance  $\bar{V}_c^m(a)$ . Specifically,

$$\frac{(\bar{R}_c^m(a) - r_c^m(a))}{\sqrt{\bar{V}_c^m(a)}} \rightarrow \mathcal{N}(0, 1) \text{ in distribution}$$

**Proof.**

Let us first rewrite the expressions in equation 6 in vector form, first introducing the following covariance matrix for  $\mathbf{p}$ :

$$\Sigma = \begin{pmatrix} Cov(p_{b_1}^{11}, p_{b_1}^{11}) & Cov(p_{b_1}^{11}, p_{b_1}^{12}) & \dots \\ Cov(p_{b_1}^{12}, p_{b_1}^{11}) & Cov(p_{b_1}^{12}, p_{b_1}^{12}) & \dots \\ \vdots & \vdots & \ddots \\ & & & Cov(p_{b_M}^{NN}, p_{b_M}^{NN}) \end{pmatrix}.$$

492 Recall that the variables  $p_a^{ij}$  are ordered lexicographically by  $(i, a, j)$ . Here we wrote  $b_1$  for the first  
 493 action in  $A$  and  $b_M$  for the last one, assuming  $|A| = M$ . Using matrix  $\Sigma$ , we can rewrite equation 6  
 494 for the approximate variance as

$$Var(R_c^m(a)) \approx (\nabla g_c^m(a)|\bar{\mathbf{p}})^T \Sigma (\nabla g_c^m(a)|\bar{\mathbf{p}}), \quad \nabla g_c^m(a)|\bar{\mathbf{p}} = \begin{bmatrix} \frac{\partial g_c^m(a)}{\partial x_{b_1}^{11}} \\ \frac{\partial g_c^m(a)}{\partial x_{b_1}^{12}} \\ \vdots \\ \frac{\partial g_c^m(a)}{\partial x_{b_M}^{NN}} \end{bmatrix} \Bigg|_{\mathbf{x}=\bar{\mathbf{p}}}, \quad (10)$$

495 where  $\nabla g_c^m(a)|\bar{\mathbf{p}}$  is the gradient vector of  $g_c^m(a)$  evaluated at  $\bar{\mathbf{p}}$ .

In the following, we employ the ‘Delta Method’ as described in [6] to allow us to derive a convergence result for the approximations for the mean and variance of  $R_c^m(a)$  that we defined above. Let us introduce a semi-vectorised representation of equation 6 where we still leverage the fact that covariances across different state-action pairs are 0, i.e.,

$$\Sigma_b^i = \begin{pmatrix} Cov(p_b^{i1}, p_b^{i1}) & Cov(p_b^{i1}, p_b^{i2}) & \dots \\ Cov(p_b^{i2}, p_b^{i1}) & Cov(p_b^{i2}, p_b^{i2}) & \dots \\ \vdots & \vdots & \ddots \\ & & & Cov(p_b^{iN}, p_b^{iN}) \end{pmatrix}$$

496 is the variance-covariance matrix for  $((p_b^{ij})_{j=1,\dots,N})$ . Since  $\Sigma$  is built by listing the  $\Sigma_b^i$  along the  
 497 diagonal for  $i = 1, \dots, N$  and  $b \in A$ , with zeros elsewhere, we have that equation 6 can be rewritten  
 498 as

$$Var(R_c^m(a)) \approx \sum_{i=1}^N \sum_{b \in A} (\nabla_b^i g_c^m(a)|\bar{\mathbf{p}})^T \Sigma (\nabla_b^i g_c^m(a)|\bar{\mathbf{p}}), \quad \nabla_b^i g_c^m(a)|\bar{\mathbf{p}} = \begin{bmatrix} \frac{\partial g_c^m(a)}{\partial x_b^{i1}} \\ \frac{\partial g_c^m(a)}{\partial x_b^{i2}} \\ \vdots \\ \frac{\partial g_c^m(a)}{\partial x_b^{iN}} \end{bmatrix} \Bigg|_{\mathbf{x}=\bar{\mathbf{p}}}, \quad (11)$$

499 where  $\nabla_b^i g_c^m(a)|\bar{\mathbf{p}}$  is the gradient vector  $(\nabla g_c^m(a)|\bar{\mathbf{p}})$  restricted to entries  $\frac{\partial g_c^m(a)}{\partial x_b^{ij}}$  for  $j = 1, \dots, N$ .

500 We refer to this approximation for the variance of  $R_c^m(a)$  as  $\bar{V}_c^m(a)$  ( $\approx Var(R_c^m(a))$ ).

Consider the random vector  $\mathbf{X} = (X_a^{ij})_{i,j=1,\dots,N}$  and  $a \in A$  (with the previously discussed lexicographic order on the  $X_a^{ij}$ ) where each  $(X_a^{ij})_{j=1}^N$  follows a Categorical distribution with probabilities  $t_a^{ij}$  - i.e.



a realisation of the vector  $\mathbf{X}$  represents the result of taking one transition from every state-action pair. Wherever  $X_a^{ij} = 1$  it represents a transition  $q^i \xrightarrow{a} q^j$ .  $\mathbf{X}$  then has means  $\mathbf{t}$  and covariances

$$\text{Cov}(X_a^{ij}, X_b^{st}) = \begin{cases} -t_a^{ij} t_b^{st} & \text{if } i = s \text{ and } a = b \\ 0 & \text{otherwise} \end{cases}$$

We can then write the variance-covariance matrix for  $\mathbf{X}$  as

$$\Sigma_{\mathbf{X}\mathbf{X}} = \begin{pmatrix} \text{Cov}(X_{b_1}^{11}, X_{b_1}^{11}) & \text{Cov}(X_{b_1}^{11}, X_{b_1}^{12}) & \dots \\ \text{Cov}(X_{b_1}^{12}, X_{b_1}^{11}) & \text{Cov}(X_{b_1}^{12}, X_{b_1}^{12}) & \dots \\ \vdots & \vdots & \ddots \\ & & & \text{Cov}(X_{b_M}^{NN}, X_{b_M}^{NN}) \end{pmatrix},$$

501 If we observe independent random samples  $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(n)}$  and denote the sample means as  
 502  $\hat{X}_b^{ij} = \frac{1}{n} \sum_{k=1}^n (X_b^{ij})^{(k)}$ , or  $\hat{\mathbf{X}} = \frac{1}{n} \sum_{k=1}^n \mathbf{X}^{(k)}$  then for the function  $g_c^m(a) [\mathbf{x}]$  we have,

$$g_c^m(a) [\hat{\mathbf{X}}] \approx g_c^m(a) [\mathbf{t}] + \sum_{i,j=1}^N \sum_{b \in A} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} (\hat{X}_b^{ij} - t_b^{ij}),$$

503 This is a direct result from the first-order Taylor expansion around  $\mathbf{t}$ , and therefore the derivatives are  
 504 evaluated at  $\mathbf{t}$ . In vector notation, we have

$$g_c^m(a) [\hat{\mathbf{X}}] \approx g_c^m(a) [\mathbf{t}] + (\nabla g_c^m(a) [\mathbf{t}])^T (\hat{\mathbf{X}} - \mathbf{t}),$$

505 where

$$(\nabla g_c^m(a) [\mathbf{t}]) = \left[ \begin{array}{c} \frac{\partial g_c^m(a)}{\partial x^{11}} \\ \frac{\partial g_c^m(a)}{\partial x_b^{12}} \\ \vdots \\ \frac{\partial g_c^m(a)}{\partial x^{NN}} \end{array} \right] \Bigg|_{\mathbf{x}=\mathbf{t}}$$

From the ‘Multivariate Delta Method’ theorem [6], as long as

$$\tau^2 := (\nabla g_c^m(a) [\mathbf{t}])^T \Sigma_{\mathbf{X}\mathbf{X}} (\nabla g_c^m(a) [\mathbf{t}]) > 0,$$

506 which we will prove later in Lemma 1 and Lemma 2, we have the following convergence:

$$\sqrt{n} \left( g_c^m(a) [\hat{\mathbf{X}}] - g_c^m(a) [\mathbf{t}] \right) \rightarrow \mathcal{N}(0, \tau^2) \text{ in distribution.} \quad (12)$$

507 Note that this is equivalent to

$$\frac{\sqrt{n} \left( g_c^m(a) [\hat{\mathbf{X}}] - g_c^m(a) [\mathbf{t}] \right)}{\tau} \rightarrow \mathcal{N}(0, 1) \text{ in distribution,} \quad (13)$$

508 where  $\tau := \sqrt{\tau^2}$ .

509 In the following we define  $\bar{\mathbf{p}}^{(n)}$  and  $\Sigma^{(n)}$  to be what  $\bar{\mathbf{p}}$  and  $\Sigma$  would have been had the agent started  
 510 with it’s prior about the transition probabilities  $\mathbf{p}$  and then witnessed exactly the transitions represented  
 511 by the random sample  $\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(n)}$ . Formally, suppose that the agent’s starting prior was, for  
 512 each state-action pair  $(q^i, b)$ , that  $p_b^{i1}, p_b^{i2}, \dots, p_b^{iN} \sim \text{Dir}(\alpha_b^{i1}, \alpha_b^{i2}, \dots, \alpha_b^{iN})$ . Then we can consider  
 513 the random variables  $p_b^{i1(n)}, p_b^{i2(n)}, \dots, p_b^{iN(n)} \sim \text{Dir}(\alpha_b^{i1} + n\hat{X}_b^{i1}, \alpha_b^{i2} + n\hat{X}_b^{i2}, \dots, \alpha_b^{iN} + n\hat{X}_b^{iN})$ .  
 514 Since  $n\hat{X}_b^{ij}$  is the count of the number of times  $X_b^{ij}$  was 1 in the random sample, this new distribution  
 515 is exactly the result of performing Bayesian inference on the prior given the random sample as our  
 516 new data. We then let

$$\bar{p}_b^{ij(n)} := \mathbb{E} \left[ p_b^{ij(n)} \right] = \frac{\alpha_b^{ij} + n\hat{X}_b^{ij}}{\sum_{k=1}^N (\alpha_b^{ik} + n\hat{X}_b^{ik})},$$

and we also define  $\Sigma^{(n)}$  as the covariance matrix of the  $p_b^{ij(n)}$  over all  $i, j, b$ , namely

$$\Sigma^{(n)} = \begin{pmatrix} Cov(p_b^{11(n)}, p_b^{11(n)}) & Cov(p_b^{11(n)}, p_b^{12(n)}) & \dots \\ Cov(p_b^{12(n)}, p_b^{11(n)}) & Cov(p_b^{12(n)}, p_b^{12(n)}) & \\ \vdots & & \ddots \\ & & & Cov(p_z^{NN(n)}, p_z^{NN(n)}) \end{pmatrix},$$

517 From Lemma 1, we have

$$\frac{\sqrt{n} \left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\hat{\mathbf{X}}] \right)}{\tau} \rightarrow 0 \text{ in probability,} \quad (14)$$

518 and this allows us to use the well-known Slutsky's Theorem [29] on equation 14 and equation 13 to  
519 show that

$$\frac{\sqrt{n} \left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\mathbf{t}] \right)}{\tau} \rightarrow \mathcal{N}(0, 1) \text{ in distribution.} \quad (15)$$

We must make one more modification to this result. Let

$$(\tau^{(n)})^2 := \left( \nabla g_c^m(a)[\bar{\mathbf{p}}^{(n)}] \right)^T \Sigma^{(n)} \left( \nabla g_c^m(a)[\bar{\mathbf{p}}^{(n)}] \right).$$

520 We would like to show that  $n(\tau^{(n)})^2 \rightarrow \tau^2$  in probability. To do this, first note that  $\bar{\mathbf{p}}^{(n)} \rightarrow \mathbf{t}$  in  
521 probability, so since  $g_c^m(a)$  has continuous derivatives we have that  $(\nabla g_c^m(a)[\bar{\mathbf{p}}^{(n)}]) \rightarrow (\nabla g_c^m(a)[\mathbf{t}])$   
522 in probability. Next we note that  $n\Sigma^{(n)} \rightarrow \Sigma_{\mathbf{X}\mathbf{X}}$  in probability. This is because for the  
523  $(i, b_1, j), (s, b_2, t)$ -entry we have  $0 \rightarrow 0$  if  $i \neq s$  or  $b_1 \neq b_2$ , and otherwise we have

$$\begin{aligned} nCov(p_b^{ij(n)}, p_b^{it(n)}) &= \frac{-n(\alpha_b^{ij} + n\hat{X}_b^{ij})(\alpha_b^{it} + n\hat{X}_b^{it})}{\left( \sum_{k=1}^N (\alpha_b^{ik} + n\hat{X}_b^{ik}) \right)^2 (1 + \sum_{k=1}^N (\alpha_b^{ik} + n\hat{X}_b^{ik}))} \\ &= \frac{-n(\alpha_b^{ij} + n\hat{X}_b^{ij})(\alpha_b^{it} + n\hat{X}_b^{it})}{\left( n + \sum_{k=1}^N \alpha_b^{ik} \right)^2 (n + 1 + \sum_{k=1}^N \alpha_b^{ik})} \\ &\rightarrow -t_b^{ij} t_b^{it} = Cov(X_b^{ij}, X_b^{it}). \end{aligned}$$

524 Therefore we have that the products converge in probability:

$$\begin{aligned} n(\tau^{(n)})^2 &= \left( \nabla g_c^m(a)[\bar{\mathbf{p}}^{(n)}] \right)^T n\Sigma^{(n)} \left( \nabla g_c^m(a)[\bar{\mathbf{p}}^{(n)}] \right) \\ &\rightarrow \left( \nabla g_c^m(a)[\mathbf{t}] \right)^T \Sigma_{\mathbf{X}\mathbf{X}} \left( \nabla g_c^m(a)[\mathbf{t}] \right) = \tau^2. \end{aligned}$$

525 Since  $\tau^2$  is always positive, and the square root function is therefore continuous at  $\tau^2$ , we have that  
526  $\sqrt{n}\tau^{(n)} \rightarrow \tau$ , and so  $\frac{\tau}{\sqrt{n}\tau^{(n)}} \rightarrow 1$  in probability. Now we can finally apply Slutsky's Theorem to  
527 obtain our final result, which is

$$\frac{\left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\mathbf{t}] \right)}{\tau^{(n)}} \rightarrow \mathcal{N}(0, 1) \text{ in distribution.} \quad (16)$$

528 Recall that  $g_c^m(a)[\mathbf{t}]$  is the actual risk in the current state  $q^c$ ,  $g_c^m(a)[\bar{\mathbf{p}}^{(n)}]$  is the agent's approximation  
529 of the expectation of the risk given it's beliefs, and  $(\tau^{(n)})^2$  is the agent's approximation of the  
530 variance of the risk given it's beliefs (both, in this case, assuming it has seen exactly  $n$  transitions  
531 from each state). So indeed our estimate of the mean of the believed risk distribution converges  
532 to the true risk with enough data, and it does so with the variance of the believed risk distribution  
533 approaching our estimate of that variance.

**Lemma 1** Given the definition of the polynomial  $g_c^m(a)[\mathbf{x}]$ , we have the following:

$$\frac{\sqrt{n} \left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\hat{\mathbf{X}}] \right)}{\tau} \rightarrow 0 \text{ in probability}$$

**Proof.**

As required for the convergence results in Theorem 3.1, one can see that all of the coefficients in  $g_c^m(a)[\mathbf{x}]$  are either 0 or 1. This means that we can rewrite it as a sum of terms of the form

$$\prod_{i,j,b} \left( x_b^{ij} \right)^{n_b^{ij}}$$

for exponents  $n_b^{ij}$ . This means that we can write

$$\frac{\sqrt{n} \left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\hat{\mathbf{X}}] \right)}{\tau}$$

534 as a sum of terms of the form

$$\frac{\sqrt{n}}{\tau} \left( \prod_{i,j,b} \left( \bar{p}_b^{ij(n)} \right)^{n_b^{ij}} - \prod_{i,j,b} \left( \hat{X}_b^{ij} \right)^{n_b^{ij}} \right).$$

535 Substituting in the definition of  $\bar{p}_b^{ij(n)}$  to this expression yields

$$\frac{\sqrt{n}}{\tau} \left( \prod_{i,j,b} \left( \frac{\alpha_b^{ij} + n \hat{X}_b^{ij}}{\sum_{k=1}^N (\alpha_b^{ik} + n \hat{X}_b^{ik})} \right)^{n_b^{ij}} - \prod_{i,j,b} \left( \hat{X}_b^{ij} \right)^{n_b^{ij}} \right)$$

536 And we can simplify this by leveraging that  $\sum_{k=1}^N (n \hat{X}_b^{ik}) = n$ , to get

$$\frac{\sqrt{n}}{\tau} \left( \prod_{i,j,b} \left( \frac{\alpha_b^{ij} + n \hat{X}_b^{ij}}{n + \sum_{k=1}^N \alpha_b^{ik}} \right)^{n_b^{ij}} - \prod_{i,j,b} \left( \hat{X}_b^{ij} \right)^{n_b^{ij}} \right)$$

537 Now, the  $\alpha_b^{ij}$  are constants, as is  $\tau$ , and the values of  $\hat{X}_b^{ij}$  are all bounded between 0 and 1. Thus to  
 538 show that this expression converges to 0 in probability, we will write it as one quotient, and show  
 539 that some term in the denominator dominates all terms in the numerator. Let  $M := \sum_{i,j,b} n_b^{ij}$ . The  
 540 expression above is equal to

$$\frac{\sqrt{n}}{\tau} \left( \frac{\prod_{i,j,b} \left( \alpha_b^{ij} + n \hat{X}_b^{ij} \right)^{n_b^{ij}} - \prod_{i,j,b} \left( \hat{X}_b^{ij} \left( n + \sum_{k=1}^N \alpha_b^{ik} \right) \right)^{n_b^{ij}}}{\prod_{i,j,b} \left( n + \sum_{k=1}^N \alpha_b^{ik} \right)^{n_b^{ij}}} \right)$$

Now on the numerator of the inner quotient, there are only two terms of order  $n^M$ . One is an

$$n^M \prod_{i,j,b} \left( \hat{X}_b^{ij} \right)^{n_b^{ij}}$$

that comes from the product on the left, and one is a

$$-n^M \prod_{i,j,b} \left( \hat{X}_b^{ij} \right)^{n_b^{ij}}$$

from the product on the right, and these cancel each other out. This means the numerator is entirely of order  $n^{M-1}$  or less. On the other hand, the denominator of the inner quotient contains the term

$n^M$ . Therefore, even after multiplying by the  $\frac{\sqrt{n}}{\tau}$  on the outside, which would mean the highest order term on in the numerator could be as high as  $n^{M-\frac{1}{2}}$ , the  $n^M$  in the denominator still dominates and the expression as a whole will converge to 0 in probability. Since

$$\frac{\sqrt{n} \left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\hat{\mathbf{X}}] \right)}{\tau}$$

541 was a sum of expressions of that form, and they all converge to 0 in probability, we get the result we  
542 desired, which is that

$$\frac{\sqrt{n} \left( g_c^m(a)[\bar{\mathbf{p}}^{(n)}] - g_c^m(a)[\hat{\mathbf{X}}] \right)}{\tau} \rightarrow 0 \text{ in probability}$$

543 **Lemma 2** The defined variable  $\tau^2 := (\nabla g_c^m(a)[\mathbf{t}])^T \Sigma_{\mathbf{X}\mathbf{X}} (\nabla g_c^m(a)[\mathbf{t}])$  is strictly greater than zero,  
544 namely  $\tau^2 > 0$ .

545 **Proof.**

546 Note that the covariance matrix can be written as  $\Sigma_{\mathbf{X}\mathbf{X}} = \mathbb{E}[(\mathbf{X} - \mathbf{t})(\mathbf{X} - \mathbf{t})^T]$  (recall  $\mathbf{t}$  is the mean  
547 vector for  $\mathbf{X}$ ). So we have

$$\begin{aligned} \tau^2 &= \mathbb{E}[(\nabla g_c^m(a)[\mathbf{t}])^T (\mathbf{X} - \mathbf{t})(\mathbf{X} - \mathbf{t})^T (\nabla g_c^m(a)[\mathbf{t}])] \\ &= \mathbb{E}[(\nabla g_c^m(a)[\mathbf{t}])^T (\mathbf{X} - \mathbf{t})^2] \end{aligned}$$

548 where we note that  $s := (\nabla g_c^m(a)[\mathbf{t}])^T (\mathbf{X} - \mathbf{t})$  is a real-valued random variable, so  $s^T = s$ . Thus  
549 to prove  $\tau^2 > 0$  we simply have to show that  $s \neq 0$  for some value of  $\mathbf{X}$  that occurs with non-zero  
550 probability.

551 Now,

$$\begin{aligned} s &= \sum_{i,j,b} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} \Big|_{\mathbf{x}=\mathbf{t}} (X_b^{ij} - t_b^{ij}) \\ &= \sum_{\text{state-action pairs } (q^i, b)} \left( \sum_{\text{possible next states } q^j} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} \Big|_{\mathbf{x}=\mathbf{t}} (X_b^{ij} - t_b^{ij}) \right) \end{aligned}$$

552 So let  $s_b^i := \sum_{\text{states } q^j} \frac{\partial g_c^m(a)}{\partial x_b^{ij}} \Big|_{\mathbf{x}=\mathbf{t}} (X_b^{ij} - t_b^{ij})$ , then  $s = \sum_{\text{state-action pairs } (q^i, b)} s_b^i$ .

553 We need to show that there is some possible value of  $\mathbf{X}$  such that  $s \neq 0$ . Now the value of  $\mathbf{X}$  is  
554 determined by the values of  $\mathbf{X}_b^i := (X_b^{ij})_{j=1}^N$  for each state-action pair  $(q^i, b)$ . Furthermore, these  $\mathbf{X}_b^i$   
555 are independent, and the value of  $s_b^i$  depends only on the value of  $\mathbf{X}_b^i$ . So if there is some state-action  
556 pair  $(q^i, b)$  such that two possible values of  $\mathbf{X}_b^i$  yield two distinct values of  $s_b^i$  both with nonzero  
557 probability, then we can fix the values of the  $X_{b'}^{hj}$  for all  $j$  and all  $(h, b') \neq (i, b)$  to be some values  
558 that occur with non-zero probability, which would fix the value of  $s - s_b^i$ , and so we could use our  
559 two distinct values of  $s_b^i$  to find two distinct values of  $s$ . Both cannot be 0, so we would be done.

560 Now, the value of  $\mathbf{X}_b^i$  is characterized by picking one  $j$  s.t.  $X_b^{ij} = 1$ , and setting all other  $X_b^{il} = 0$  for  
561  $l \neq j$ . This means that to find two different values of some  $s_b^i$ , we just need to find states  $q^i, q^j, q^l$  and  
562 an action  $b$  such that the derivatives  $\frac{\partial g_c^m(a)}{\partial x_b^{ij}} \Big|_{\mathbf{x}=\mathbf{t}}$  and  $\frac{\partial g_c^m(a)}{\partial x_b^{il}} \Big|_{\mathbf{x}=\mathbf{t}}$  are distinct. Then setting  $X_b^{ij} = 1$   
563 would yield a different value of  $s_b^i$  from setting  $X_b^{il} = 1$ . So long as the events  $X_b^{ij} = 1$  and  $X_b^{il} = 1$   
564 both have nonzero probability, we would be done.

In order to show that such states  $q^i, q^j, q^l$  and such an action  $b$  exist, we must introduce vectors  $A^n$   
that will effectively keep track of each state's contribution towards  $g_c^m(a)[\mathbf{t}]$  at the  $n$ th step of the  
risk backpropagation. First, define the  $N$ -by- $N$  matrix  $P'_n[\mathbf{x}]$  for  $n = 0, 1, \dots, m-2$  such that

$$(P'_n[\mathbf{x}])_{ij} = \begin{cases} 1 & \text{if } i = j \text{ and } q^i \text{ is unsafe and observed} \\ 0 & \text{if } i \neq j \text{ and } q^i \text{ is unsafe and observed} \\ x_{b_i n}^{ij} & \text{otherwise} \end{cases}$$

565 where where  $b_{in} := \arg \min_b \bar{R}_i^n(b)$ . Define  $P'_{m-1}[\mathbf{x}]$  as

$$(P'_{m-1}[\mathbf{x}])_{ij} = \begin{cases} 1 & \text{if } i = j \text{ and } q^i \text{ is unsafe and observed} \\ 0 & \text{if } i \neq j \text{ and } q^i \text{ is unsafe and observed} \\ x_a^{ij} & \text{otherwise} \end{cases}$$

566 Then the  $P'_n[\mathbf{x}]$  represent the transition probabilities used in the calculation of  $g_c^m(a)[\mathbf{x}]$ . Specifically,  
567 we have that

- 568 •  $g_k^n[\mathbf{x}]$  is the  $k$ th entry of the vector  $(P'_{n-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}])g^0$  for  $n < m$
- 569 •  $g_k^m(a)[\mathbf{x}]$  is the  $k$ th entry of the vector  $(P'_{m-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}])g^0$
- 570 • So the risk at current state  $q^c$ ,  $g_c^m(a)[\mathbf{t}]$ , is the  $c$ th entry of the vector
- 571  $(P'_{m-1}[\mathbf{t}]) \dots (P'_0[\mathbf{t}])g^0$

572 where  $g^0$  is the vector with entries  $(g^0)_k := \mathbb{1}(q^k \text{ is observed and unsafe})$ . We can now define the  
573 vectors  $A^n$  for  $n \leq m$  by

$$A_i^n := \begin{cases} ((P'_{n-1}[\mathbf{t}]) \dots (P'_0[\mathbf{t}])g^0)_i & \text{if } q^i \text{ is safely reachable from } q^c \text{ in exactly } m - n \text{ steps} \\ 0 & \text{otherwise} \end{cases}$$

574 Where in this case a state  $q^{s_n}$  is defined to be *safely reachable* from the current state  $q^{s_0} = q^c$  in  
575 exactly  $n$  steps if

- 576 • there are states  $q^{s_1}, q^{s_2}, \dots, q^{s_{n-1}}$  such that each  $t_{b_{s_1}}^{s_p s_{p+1}} > 0$  for actions  $b_{s_0} = a$  and
- 577  $b_{s_k} := \arg \min_b \bar{R}_{s_p}^{m-k-1}(b)$  determined by the agent's expected safest policy, and
- 578 • the states  $q^{s_1}, q^{s_2}, \dots, q^{s_{n-1}}$  are all safe (note that  $q^{s_n}$  can still be unsafe)

579 The purpose of these  $A^n$  is just to restrict our attention to the states at step  $n$  of the backpropagation  
580 that actually influence  $g_c^m(a)[\mathbf{t}]$ . It is easy to see that

$$((P'_{m-1}[\mathbf{t}]) \dots (P'_n[\mathbf{t}])A^n)_c = g_c^m(a)[\mathbf{t}] \text{ for every } n = 0, 1, \dots, m \quad (17)$$

581 Now we will be able to argue that if  $g_c^m(a)[\mathbf{t}]$  is not equal to 0 or 1, there are states  $q^i, q^j, q^l$  and  
582 an action  $b$  such that  $t_b^{ij}$  and  $t_b^{il}$  are both non-zero (so there is a positive probability of the events  
583  $X_b^{ij} = 1$  and  $X_b^{il} = 1$ ) and such that  $\left. \frac{\partial g_c^m(a)}{\partial x_b^{ij}} \right|_{\mathbf{x}=\mathbf{t}} > \left. \frac{\partial g_c^m(a)}{\partial x_b^{il}} \right|_{\mathbf{x}=\mathbf{t}}$ .

584 So assume that  $g_c^m(a)[\mathbf{t}]$  is not equal to 0 or 1. Let  $n_0$  be the largest index such that  $A^{n_0}$  contains an  
585 entry  $(A^{n_0})_l$  that is equal to 0 and such that  $q^l$  is safely reachable from  $q^c$  in exactly  $m - n_0$  steps -  
586 so  $(A^{n_0})_l$  is a 0 that came from  $(P'_{m-1}[\mathbf{t}]) \dots (P'_0[\mathbf{t}])g^0$ .

587 Since  $g_c^m(a)[\mathbf{t}]$  is not 0,  $n_0 < m$ , and since  $q^l$  is safely reachable in  $m - n_0$  steps, let  $q^c =$   
588  $q^{s_0}, q^{s_1}, \dots, q^{s_{m-n_0}} = q^l$  be a path along which  $q^l$  is safely reachable. Then let  $q^i = q^{s_{m-n_0-1}}$ , and  
589 we have that  $q^i$  is safe, and  $t_{b_{s_{m-n_0-1}}}^{il} > 0$ . For brevity, write  $b' := b_{s_{m-n_0-1}}$

590 Now since  $q^i$  is safely reachable in  $m - (n_0 + 1)$  steps,  $(A^{n_0+1})_i$  cannot be equal to 0 (since  $n_0$   
591 was maximal), so there must be some state  $q^j$  such that  $t_{b'}^{ij} > 0$  and  $A_j^{n_0} > 0$ , (in order for the term  
592  $t_{b'}^{ij} A_j^{n_0}$  to contribute some positive value to  $A_i^{n_0+1}$ ). Finally, let  $p$  be the probability of safely entering  
593  $q^i$  in  $m - (n_0 + 1)$  steps (i.e., the sum over all paths that safely reach  $q^i$  of the probability of taking  
594 that path by choosing the actions specified by the agent's expected safest policy). Then by the chain  
595 rule,

$$\left. \frac{\partial g_c^m(a)}{\partial x_{b'}^{ij}} \right|_{\mathbf{x}=\mathbf{t}} = p \left( 1 \times A_j^{n_0} + t_{b'}^{ij} \times \left. \frac{((P'_{n_0-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}])g^0)_j}{\partial x_{ij}} \right|_{\mathbf{x}=\mathbf{t}} \right) > 0$$

596 since clearly  $\left. \frac{((P'_{n_0-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}])g^0)_j}{\partial x_{ij}} \right|_{\mathbf{x}=\mathbf{t}}$  cannot be negative. On the other hand,

$$\begin{aligned} \left. \frac{\partial g_c^m(a)}{\partial x_{b'}^{il}} \right|_{\mathbf{x}=\mathbf{t}} &= p \left( 1 \times (A^{n_0})_l + t_{b'}^{il} \times \left. \frac{((P'_{n_0-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}]) g^0)_l}{\partial x_{b'}^{il}} \right|_{\mathbf{x}=\mathbf{t}} \right) \\ &= p \left( 1 \times 0 + t_{b'}^{il} \times \left. \frac{((P'_{n_0-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}]) g^0)_l}{\partial x_{b'}^{il}} \right|_{\mathbf{x}=\mathbf{t}} \right) = 0 \end{aligned}$$

597 since only one of  $t_{b'}^{il}$  and  $\left. \frac{((P'_{n_0-1}[\mathbf{x}]) \dots (P'_0[\mathbf{x}]) g^0)_l}{\partial x_{b'}^{il}} \right|_{\mathbf{x}=\mathbf{t}}$  can be nonzero - if increasing the value of  $t_{b'}^{il}$   
 598 could increase the value of  $(A^{n_0})_l = ((P'_{n_0-1}[\mathbf{t}]) \dots (P'_0[\mathbf{t}]) g^0)_l$  from 0 to greater than 0, then  $t_{b'}^{il}$   
 599 must have been 0 since  $((P'_{n_0-1}[\mathbf{t}]) \dots (P'_0[\mathbf{t}]) g^0)_l$  is a sum of products of values from  $\mathbf{t}$ , all of which  
 600 are non-negative.

601 Hence we have found states  $q^i, q^j, q^l$  and an action  $b'$  such that the derivatives  $\left. \frac{\partial g_c^m(a)}{\partial x_{b'}^{ij}} \right|_{\mathbf{x}=\mathbf{t}}$  and  
 602  $\left. \frac{\partial g_c^m(a)}{\partial x_{b'}^{il}} \right|_{\mathbf{x}=\mathbf{t}}$  are distinct. Hence the claim.

603 The only detail left to note is that we assumed that  $g_c^m(a)[\mathbf{t}]$  is not either equal to 0 or 1. This  
 604 assumption is reasonable to make, because if it did not hold, then either our agent would be doomed  
 605 to enter an unsafe state within  $m$  steps, or there is no chance of entering an unsafe state within  $m$   
 606 steps, according to the agent's expected safest actions. Since what matters to us is how the agent  
 607 manages risk, situations involving risk 1 or risk 0 are irrelevant.

608 **F Appendix F. Confidence Bound on the Risk**

609 To estimate a confidence bound on the risk, we appeal to the Cantelli Inequality, which is a one-sided  
610 Chebychev bound [5], and states that for a real-valued random variable  $R$  with expectation  $\mathbb{E}[R]$  and  
611 variance  $\text{Var}[R]$ , for  $\lambda > 0$  we have

$$\Pr(R \leq \mathbb{E}[R] + \lambda) \geq 1 - \frac{\text{Var}[R]}{\text{Var}[R] + \lambda^2}$$

612 If we let  $C := 1 - \frac{\text{Var}[R]}{\text{Var}[R] + \lambda^2}$ , then rearranging we get that  $\lambda = \sqrt{\frac{\text{Var}[R]C}{1-C}}$ . Thus for a variable  $R$   
613 that represents some sort of risk, and for some value of  $0 < C < 1$ , we can say

$$\Pr(R \leq P) \geq C$$

614 where  $P := \mathbb{E}[R] + \sqrt{\frac{\text{Var}[R]C}{1-C}}$ . In words, “there is at least  $C$  chance that the risk is at most  $P$ .”  
615 Alternatively, “we are at least  $\frac{C}{100}\%$  confident that the risk is at most  $P$ .”

616 **G Appendix G**

617 To understand what exactly  $\bar{R}_c^m(a)$  is an approximation of, consider instead calculating this risk  
 618 using the true transition probabilities  $t_a^{kj}$ , We would get

$$r_k^0 := \mathbb{1}(q^k \text{ is observed and unsafe}) \tag{18}$$

$$r_k^{n+1}(a) := \begin{cases} 1 & \text{if } q^k \text{ is observed and unsafe} \\ \sum_{j=1}^N t_a^{kj} r_j^n & \text{otherwise} \end{cases} \tag{19}$$

$$r_k^{n+1} := \begin{cases} 1 & \text{if } q^k \text{ is observed and unsafe} \\ r_k^{n+1}(\arg \min_{a \in A} \bar{R}_k^{n+1}(a)) & \text{otherwise} \end{cases} \tag{20}$$

619 Note that we crucially still take the minimum risk action  $a$  according to the agent’s approximation  
 620  $\bar{R}_k^{n+1}(a)$ . In this case, the term  $r_c^m(a)$  is the true probability of entering an unsafe state after selecting  
 621 action  $a$  in the agent’s current state  $q^c$  and thereafter selecting the actions that *the agent currently*  
 622 *believes* will minimize the probability of entering an unsafe state over the horizon  $m$ .  $\bar{R}_c^m(a)$  is the  
 623 agent’s approximation of  $r_c^m(a)$ .

624 We will later justify the use of  $\bar{R}_c^m(a)$  as an approximation of  $r_c^m(a)$ , but for now let us consider why  
 625 it makes sense to define  $m$ -step risk as  $r_c^m(a)$ . This because the action  $a$  that minimizes believed risk  
 626 is the action that the agent would choose if it was trying to behave as safely as possible, what I will  
 627 call going into ‘safety mode’. Consider the motivating example of a pilot learning to fly a remote  
 628 control helicopter by incrementally expanding the set of actions they feels safe taking. They start by  
 629 generating just enough lift to begin flying, then immediately land back down again. They repeat this  
 630 process a few times until they feel that they have a good understanding of how the helicopter responds  
 631 to this limited range of inputs. Then they take a risk (by either flying a bit higher, or attempting to  
 632 move horizontally) and once again immediately land. As they repeat this process of taking small risks  
 633 and landing to remain safe, they begin to expand their comfort zone. At some point after taking a  
 634 risk, they will feel comfortable just coming back to a hovering position rather than landing, once they  
 635 have become confident that they can hover safely. This suggests that a natural process for learning  
 636 to operate in the face of risks is to repeatedly take small risks followed by going into safety mode  
 637 until back in a confidently safe state. Thus, when calculating how risky an action is, it makes sense  
 638 to consider the probability of entering an unsafe state given that after the action the agent will enter  
 639 safety mode.  $r_c^m(a)$  does exactly this.

640 As mentioned earlier, the other reason for defining the risk  $r_c^m(a)$  in this way is that it makes it possible  
 641 for the agent to attempt to calculate the risk without having to reason about the inter-dependency  
 642 between the calculated risk and the agent’s future actions. However, it does more than this. We  
 643 will see in the next section that it in fact allows the agent to view  $\bar{R}_c^m(a)$  as (an approximation of)  
 644 the expected value of a random variable for the believed risk, where we can also approximate the  
 645 *variance* of that random variable, allowing for deeper reasoning about action-selection for Safe RL.