

Optimistic Risk-Aware Model-based Reinforcement Learning

Romina Abachi *University of Toronto, Vector Institute*

romina@cs.toronto.edu

Amir-massoud Farahmand *University of Toronto, Vector Institute*

farahmand@vectorinstitute.ai

Abstract

We study the problem of planning with a learned model in model-based RL. In this setting, the agent should act safely with respect to its belief about the environment, as it may have errors in that belief. On the other hand, it needs to explore in order to learn good policies. We propose two algorithms based on the notion of risk-sensitivity over the uncertainty over parameters of the true MDP using a Bayesian RL framework. In the first, we propose optimizing the risk over the beliefs directly using a policy sub-gradient. In the second, we propose using the risk as a constraint to an optimistic or exploratory planner. We empirically evaluate our proposed algorithms and see improvement in certain cases over non-risk-aware and only-optimistic methods.

Keywords: Model-based Reinforcement Learning, Risk, Model-based, Exploration-Exploitation

1. Introduction

Model-based reinforcement learning (MBRL) consists of learning a model of the environment, planning using the learned model, and using the obtained policy to further interact with the environment, collect more data, and further improve the model and policy (Sutton, 1990). Variations of this procedure have gained considerable interest in recent years (Schrittwieser et al., 2020; Hafner et al., 2019; Ha and Schmidhuber, 2018).

A learned model, however, is often imprecise – it has some errors compared to the true dynamics of the environment. The sources of errors are the estimation error (i.e., the statistical error caused by having a finite number of samples) and possibly model approximation error (i.e., the true dynamics cannot be represented by any member of the model class due to limited capacity), see e.g., Györfi et al. (2002); Steinwart and Christmann (2008) for a discussion of estimation and function approximation errors in the supervised learning setting, Munos and Szepesvári (2008); Antos et al. (2008); Farahmand et al. (2016); Tosatto et al. (2017); Chen and Jiang (2019) in model-free RL, and Farahmand et al. (2017) in MBRL. Many MBRL algorithms ignore this error, and plan as if the learned model is the true model of the environment. This approach is called the *certainty equivalence principle* in control theory (Aström and Wittenmark, 1994). The underlying assumption of the certainty equivalence principle is that if and when the model becomes accurate enough, the difference between the policy that is optimal according to the model versus the true optimal policy of the environment becomes negligible. If there is no model approximation error, this might be asymptotically true. But oftentimes we cannot guarantee that there is no model approximation error when we model a complex environment. Moreover, in the finite sample regime, the estimation error might still be substantial. As a consequence, following this recipe may lead to a policy whose behaviour is disastrous in the real environment. The planning algorithm may exploit the errors in the model to improve its performance on the model, but these errors could lead to bad and unsafe outcomes in the true environment.

To mitigate this issue, we can either improve the model learning or use planners that explicitly incorporate the error in the model. To improve the model learning, we can use models that are more expressive (e.g., larger deep neural networks) or use decision-aware model learning approaches that focus on learning the model where it matters for the decision making while ignoring irrelevant aspects of the environment (Farahmand et al., 2017; D’Oro et al., 2020; Schrittwieser et al., 2020; Grimm et al., 2020; Lambert et al., 2020; Eysenbach et al., 2021). Although this is a worthy approach, in this work we instead turn our attention to the other part of the problem: planning with an imperfect model.

To deal with an imperfect model, we propose and study an online Bayesian risk-sensitive method. The method maintains a posterior over possible models of the environment, as is common in Bayesian RL, and the planner finds a policy that is risk-sensitive with respect to the uncertainty over the true model. Specifically, we use Conditional Value-at-Risk (CVaR) as the risk measure (Acerbi and Tasche, 2002; Shapiro et al., 2009). As opposed to the performance

measured according to the expected value function according to the posterior over the models, which is common in the Bayesian RL framework, CVaR measures the average performance of only those models whose value function belongs to a lower quantile of performance according to the posterior. This allows control over how risk-averse the policy is. Our first algorithmic contribution is a policy gradient-based algorithm for risk-sensitive planning (Section 3). This risk-aware Bayesian RL viewpoint has been studied in some previous work (cf. Angelotti et al. (2021) and Chapter 6 of Ghavamzadeh et al. 2015), but the specific PG-based method has not been proposed before (among some other differences). This approach has also some high-level similarities with the Robust Dynamic Programming (Iyengar, 2005; Nilim and El Ghaoui, 2005) literature, which focuses on robustifying against the worst-case statistically feasible model. The usual notion of robustness, however, can be too conservative as it considers the worst-case model. CVaR, on the other hand, provides a gentler notion of risk that considers a whole spectrum between the expected performance (risk neutral) to almost worst-case. We also note that our usage of risk is different from what is common in risk-sensitive RL, which focuses on uncertainty of return within a fixed MDP, see e.g., Chow et al. (2017) and references therein. In our work, the uncertainty is over the true model of the environment.

To allow the MBRL agent to explore the environment while avoiding risky behaviours, we introduce Optimistic Risk-Constrained planner in Section 4 as the other algorithmic contribution of this work. It modifies the optimistic planners and its stochastic variants such as PSRL (Osband and Van Roy, 2017), which are commonly used for balancing exploration-exploitation, by explicitly considering risk over uncertainty of the model within optimistic planning. We empirically evaluate a few variants and compare them with purely optimistic and purely risk-averse methods. We find that a CVaR constraint can improve performance in certain domains while having little effect in others.

2. Preliminaries

2.1 Markov Decision Processes

We consider the problem of learning to optimize a discounted infinite-horizon MDP $(\mathcal{X}; \mathcal{A}; \gamma; \mathcal{P}; R; \mu_0)$ with state space \mathcal{X} , action space \mathcal{A} , discount factor γ , and starting state distribution μ_0 (Bertsekas and Tsitsiklis, 1996; Szepesvári, 2010; Sutton and Barto, 2018). The transition kernel $P: \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{X})$ and reward function $R: \mathcal{X} \times \mathcal{A} \times \mathcal{X} \rightarrow \mathbb{R}$ are unknown, where $\mathcal{P}(\cdot)$ refers to the space of probability distributions defined over a set \cdot . The RL agent interacts with the MDP starting at time 0, takes action A_t at each timestep according to a policy π , transitions to state X_{t+1} , and observes a scalar reward $R_t = R(X_t, A_t, X_{t+1})$. We want to find a stationary Markov policy $\pi_M: \mathcal{X} \rightarrow \mathcal{P}(\mathcal{A})$ that maximizes policy performance $J_0(\pi; P) = \sum_{x \in \mathcal{X}} \mu_0(x) V_\pi(x) = \mathbb{E}_{P; \mu_0} \sum_{t=0}^{\infty} \gamma^t R(X_t, A_t, X_{t+1})$; where $V_\pi(x)$ is defined as the value function of the MDP with transition kernel P at state x . In the rest of the paper we omit μ_0 in our notation for J_0 . If there is no restriction over the policy, the optimal policy does not depend on the initial distribution μ_0 . On the other hand, if the policy is restricted to a subset of all possible policies, the optimal policy within that class may depend on μ_0 . If the policy is parameterized by θ , we simply use the parameters of the policy to refer to the policy, i.e., the policy performance and value functions would be $J_\theta(\pi; P)$ and $V_\pi(\cdot)$, respectively.

Suppose we have collected all the interactions of the environment in a dataset consisting of tuples $(X_t; A_t; R_t; X_{t+1})$. We consider the Bayesian setup where given a prior distribution over unknown MDPs (transition kernel and reward function), we infer the posterior distribution $(M|D)$. The posterior represents our belief about the true MDP given the collected data so far. As more data is collected, the posterior concentrates more and more around μ_0 of the true MDP, assuming that prior probability is non-zero, i.e., $\mu_0(M) > 0$. In much of this work we omit the dependence of the posterior and policy performance on the unknown function μ_0 only for notational convenience. However, we assume a posterior over reward functions is also inferred.

2.2 Risk Measures

A risk measure maps a random variable to a real-valued number. It is a measure of how risky a random variable behaves. Various risk measures have been proposed in the literature, see e.g., Shapiro et al. (2009). In this work, we focus on the Value-at-Risk (VaR) and Conditional Value-at-Risk (CVaR).

Let Z be a scalar random variable defined on a probability space $(\mathcal{F}; \mathcal{P})$ with a finite mean, i.e., $E[Z] < \infty$. Denote its cumulative distribution function (CDF) by $F_Z(z) = P(Z \leq z)$. The Value-at-Risk (VaR) with confidence level $\alpha \in (0; 1)$ is defined as: $\text{VaR}_\alpha(Z) = \inf\{z \mid F_Z(z) \geq \alpha\}$. This is the greatest point $z \in \mathbb{R}$ at which the probability of random variable Z taking values smaller than z is $\leq 1 - \alpha$. This quantity is sometimes also referred to as the (lower) α -quantile. Note that as α increases, $\text{VaR}_\alpha(Z)$ is non-decreasing and right-continuous, and the maximum is attained.

To define the Conditional Value-at-Risk, we define an auxiliary functional first:

$$F(Z; t) = \int_{-\infty}^t F_Z(z) dz + \frac{1}{\alpha} \int_t^{\infty} F_Z(z) dz; \quad (1)$$

where $\int_{-\infty}^t f(z) dz = \min\{0, t\}$ (we define $\int_t^{\infty} f(z) dz = \max\{0, t\}$). The minimizer of this functional is CVaR of the random variable at level α , as shown by [Rockafellar et al. \(2000\)](#):

$$\text{CVaR}_\alpha(Z) = \min_{z \in \mathbb{R}} F(Z; z); \quad (2)$$

The minimum is attained on a non-empty closed bounded interval, with $\text{VaR}_\alpha(Z)$ on the left endpoint of that interval (Theorem 1 of [Rockafellar et al. \(2000\)](#) or Section 6.2.4 of [Shapiro et al. 2009](#)). Therefore, we have $\text{CVaR}_\alpha(Z) \geq \text{VaR}_\alpha(Z)$. When the CDF $F_Z(z)$ is continuous at $\text{VaR}_\alpha(Z)$, which means that there is no probability atom there, CVaR takes a simpler form of

$$\text{CVaR}_\alpha(Z) = E[Z \mid Z \geq \text{VaR}_\alpha(Z)]; \quad (3)$$

see [Acerbi and Tasche \(2002\)](#) or Theorem 6.2 of [Shapiro et al. \(2009\)](#). Thus, $\text{CVaR}_\alpha(Z)$ can be interpreted as the expected value of Z up to $\text{VaR}_\alpha(Z)$ on the lower α -quantile of the distribution. Note that $\text{CVaR}_0(Z) = E[Z]$, and $\text{CVaR}_\alpha(Z)$ tends to $\text{ess\,inf}(Z)$ as $\alpha \rightarrow 1$. CVaR is popular in financial applications in comparison to the VaR as it is useful for controlling rare, but potentially catastrophic events, that occur below the α -quantile and do not impact the VaR at all ([Rockafellar et al., 2006](#)). As opposed to VaR, CVaR is a coherent risk measure which gives it additional desirable properties ([Artzner et al., 1999](#)). As a result, we focus on optimizing CVaR in this work.

3. Risk-sensitive Planning for MBRL

We describe the risk-aware planning for model-based RL and propose our policy gradient-based algorithm for solving it. The risk in this context refers to parametric (epistemic) uncertainty: the uncertainty due to the unknown parameters of the MDP. This is in contrast to the inherent uncertainty present when interacting with a known MDP, that is, the uncertainty due to the return distribution, induced by the stochasticity of the environment and the policy (aleatoric uncertainty). This framework is sometimes referred to as Risk-aware Bayesian RL ([Ghavamzadeh et al., 2015](#), Section 6).

Figure 1 illustrates why it is important to consider a risk-sensitive objective. Suppose that the agent has already collected some data through interaction with the true MDP with the dynamics \mathcal{M} . In a Bayesian framework, the agent's belief about the probability of each model in the space of all models is encoded through its posterior $(P_j|D)$. This uncertainty over the model induces an uncertainty on the performance of a policy π for any $P_j \in (P_j|D)$, the agent evaluates the performance $J(\pi; P)$. This defines a distribution $(J(\pi) | D)$ over the performance of this policy, as indicated by a green curve in the figure.

Now suppose that the agent wants to select one of the two policies π_1 and π_2 . Which one should it choose? One approach is to compare the average performance w.r.t. the posterior. This would be $E_{(P_j|D)}[J(\pi_1; P)]$ and $E_{(P_j|D)}[J(\pi_2; P)]$. In this figure, π_2 would be the preferred policy according to this measure. Looking at the distribution, however, we see that $(J(\pi_2) | D)$ has a fatter tail in the lower values of performance. This means that according to our posterior belief, it is more probable that policy π_2 's performance would be significantly worse than π_1 's, even though its average performance is not. If we define "safe" as the policy with a lower chance of low performances across the model space (and "risky" as vice versa), policy π_1 is safer than π_2 according to our posterior belief. This discussion leads us to the definition of a risk-sensitive objective and policy.

Definition 1 Given a posterior distribution $(P_j|D)$ over the model of the MDP, starting state distribution ρ and risk measure ρ , the risk-sensitive objective is defined as

$$J_\rho(\pi; P) = \int_{-\infty}^{\infty} \rho(z) J(\pi; P) dz; \quad (4)$$

Figure 1: The distribution of policy performance for two policies π_1 and π_2 induced by the the posterior over models given D_n . Policy π_2 has higher expected performance, but higher probability of low-performance over models, i.e., it is less safe. On the right side of the figure we also illustrate which parts of the posterior are considered for the optimistic algorithms defined in Section 4.

The risk-sensitive policy over the policy space is defined as

$$\arg\max_{\pi} \mathbb{E}_{\mathbb{P}} [J(\pi; \mathbb{P})] \tag{5}$$

Usually, \mathbb{E} is simply taken as the expectation (Zintgraf et al., 2019), which is a risk-neutral measure. In this work we are interested in risk measures that lead to risk-averse behaviour. In particular, we choose CVaR, as introduced in Section 2.2, which is a versatile and commonly used risk measure. In our setting, the random variable $Z = J(\pi; \mathbb{P})$, the policy performance under policy π and transition kernel $\mathbb{P}(\cdot | \mathcal{D})$.

Our proposed risk-sensitive planning for MBRL algorithm, which is specified in detail in Algorithm 1, at the high-level works as follows: The agent interacts with the environment, and maintain a posterior over the model and the reward function. It then computes the (sub-) gradient of CVaR($J(\pi)$) w.r.t. θ in order to update the policy. This procedure repeats. We explain the details in the rest of this section.

Computing CVaR($J(\pi)$) exactly, which involves computing the integral (1), is not feasible. Instead, we use an empirical estimate of (1) using m sampled models $\mathbb{P}_i(\cdot | \mathcal{D})$ ($i = 1; \dots; m$):

$$\hat{F}(\pi; \theta) = \mathbb{E}_{\mathbb{P}} [J(\pi; \mathbb{P})] + \frac{1}{m} \sum_{i=1}^m J(\pi; \mathbb{P}_i) \tag{6}$$

We can compute the empirical estimate of CVaR($J(\pi)$) by noting that it is equal to $\hat{F}(\pi; \theta)$ with $\theta = \text{VaR}(J(\pi))$. The value of $\text{VaR}(J(\pi))$ itself is estimated as follow: we compute the performance $J(\pi; \mathbb{P}_i)$ ($i = 1; \dots; m$), sort them in increasing order $J_{(1)} \leq J_{(2)} \leq \dots \leq J_{(m)}$, and then estimate $\text{VaR}(J(\pi))$ as $J_{(bm)}$. The empirical estimate of CVaR($J(\pi)$) is then equal to $\text{mean}(J_{(1)}; \dots; J_{(bm)})$.

To perform the policy (sub-) gradient over CVaR($J(\pi)$), we need to compute the sub-gradient $\partial F(\pi; \theta)$ with respect to θ . We use the property that the posterior $\mathbb{P}(\cdot | \mathcal{D})$ is not dependent on θ . Let $\hat{\theta}$ be an estimate of VaR as specified above. Then the sub-gradient of (6) with respect to θ is given by:

$$\partial \hat{F}(\pi; \theta) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{J(\pi; \mathbb{P}_i) \leq \theta\}} \tag{7}$$

In the rest of the paper, we may overload the notation $\mathbb{1}_{\{J(\pi; \mathbb{P}_i) \leq \theta\}}$ to refer to the components that are being compared to determine its value. Note that this sub-gradient computation is simpler than the procedure introduced by Tamar et al. (2015), in which the distribution also depends on the parameter from which we take the gradient. Our procedure is similar to the procedure used by Rockafellar et al. (2000) in their application to hedging.

How many samples should be collected? We want to have enough samples for the estimates to be accurate while not wasting storage and computation on collecting and solving (i.e. through policy evaluation) more samples than necessary. To answer this question, [Briggs and Ying \(2018\)](#) propose a stopping condition based on distribution-free confidence intervals. This is the condition we use on L24 and L25 of Algorithm 1. The relative error tolerance controls how wide the confidence intervals are determined with what confidence level the quantile falls in the confidence interval. For example, if we want 95% confidence intervals, we would set $\epsilon = 0.05$, and depending on ϵ_{rel} , we would need to collect enough data to get narrow enough intervals and 95% confidence that the ϵ -quantile falls in those intervals. We note that using this stopping condition to estimate risk measures in RL is similar to how it is used by [Angelotti et al. \(2021\)](#), which we discuss in more depth at the end of this section.

Our overall approach is presented in Algorithm 1. We use the notation μ_j to refer to the j th item in the sequence V . In each iteration of the algorithm, we collect some transitions by acting on the environment with the current policy and update the posteriors $\mathcal{P}; \mathcal{R}$ (L4 - L5). We then take policy (sub-)gradient (PG) steps to optimize the objective (4) (L7-L11), starting from the previous iteration's best policy $\mu_{(best)}^{(k)}$ as the initial point for optimization in the next round. We use the function EVALUATE RISK in each step of PG to evaluate the risk for each set of policy parameters. This function also calculates and returns the PG for each of the posterior samples used to estimate the (C)VaR. $r_{V_k}[i]$ is the PG of the i th value function in the sorted V , which may be computed using any PG method such as REINFORCE or an actor-critic algorithm. Note that although in each call to EVALUATE RISK, new samples from the posterior are collected, we found empirically that collecting a large number of samples from the posterior after each update of the posteriors and performing the PG step on those samples was much more stable for optimization. Note that after each update of the posterior, the risk-sensitive objective (4) and the corresponding optimizer (5) change. Since the subgradient is not necessarily in the direction of increasing risk, we keep track of the best iterate $\mu_{best}^{(k)}$ among $\mu_1; \dots; \mu_k$, i.e., $\mu_{(best)}^{(k)} = \max_{j=1; \dots; k} \mu_j$; where μ_j refers to the VaR estimate of μ_j .

Solving the optimization problem (5) is also considered by [Angelotti et al. \(2021\)](#). Similar to our work, the approximate the risk measure by sampling from the posterior. There are, however, a few key differences. Their problem setup is the ofline (batch) setting, and not the online setting of this work. They do not empirically study the effect of finding a risk-sensitive policy on the performance of a continually interacting agent, as we do later in Section 5. From the algorithmic perspective, their approach to finding a good policy is based on an exhaustive search in the space of all policies or a subset thereof. On the other hand, the algorithm in this section is a policy gradient-based approach that gradually updates the policy in the direction of improvement.

4. Optimistic Risk-Constrained Planning for MBRL

There is a tradeoff between performance and safety. An agent that behaves safely with respect to its belief about the environment may be hesitant to learn much about its environment. This causes it to act too conservatively and degrades its performance. For example, consider a navigating robot in an unknown terrain. If the robot encounters a crevice or hole in the ground on its first few steps, its model of the environment will be biased towards assuming that it is surrounded by unsafe terrain. A risk-averse planner may lead to a policy that causes it to simply stay in one spot, and never accomplish its intended task.

To have a high performance, the agent needs to carefully balance its exploration and exploitation. This can be done through optimistic methods such as UCRL methods ([Jaksch et al., 2010](#)) or Thompson/Posterior sampling ([Strens, 2000](#); [Osband et al., 2013](#)), which can be interpreted as a stochastic optimistic approach ([Osband and Van Roy, 2017](#)). Nonetheless, this may lead to some behaviours that are risky according to the agent's belief about its environment.

Our desire is to have an agent that has good performance while reducing the number of catastrophic events that it faces. It is thus natural to consider safety as a constraint rather than its own objective. We formulate this as follows:

$$\begin{aligned} \text{OptimisticPlanner}(\mu; J) & \quad (8a) \\ \text{subject to } (J(\mu; P)) & \leq b; \quad (8b) \end{aligned}$$

where $\text{OptimisticPlanner}(\mu; J)$ is an optimistic planner, J is a risk measure, and b is the acceptable level of risk. The optimistic planner depends on the posterior $(\mathcal{P}; \mathcal{D})$ and the policy performance $(\mu; P)$ and may be formulated in different ways. The risk measure can be replaced with any risk measure over the posterior $(\mathcal{P}; \mathcal{D})$. Here, we use the CVaR as in the previous section.

Algorithm 1 PG Algorithm for Risk-sensitive Planning

```

1: Input: prior distributions  $P_0(P)$ ;  $R_0(R)$ , quantile order  $\alpha \in (0; 1)$ , relative error tolerance  $\epsilon_{rel}$ , initial state
   distribution  $\mu_0$ , confidence level  $\beta$ , policy learning rates  $(\eta)$ , convergence criterion  $\delta$ , initial parameters  $\theta_1$ .
2: Initialization: policy parameters  $\theta_{(best)}^1 = \theta_1$ ,  $\text{Dataset} = \text{fg}$ .
3: for  $k = 1; 2; \dots$  do
4:   Collect transition(s) with policy  $\pi_k$  from the environment and add to  $\text{data} := \mathcal{D} [f(X_t; A_t; X_{t+1}); R_t]$ .
5:   Update posteriors  $P^k(P|D)$ ;  $R^k(R|D)$ .
6:   Policy parameters  $\theta_k = \theta_{(best)}^k$ . Re-initialize  $\epsilon_k = 1$ .
7:   repeat
8:      $V_k; r = V_k; \pi_k; \hat{F}^{\alpha}(k; \pi_k) = \text{EVALUATE RISK}(k; \epsilon_k; \beta; P^k; R^k)$ 
9:      $\pi_{k+1} = \pi_k + (k) \frac{1}{N_k} \sum_{i=1}^{N_k} (V_k[i]; \pi_k) r - V_k[i]$ 
10:     $k = k + 1$ 
11:   until  $k - \pi_k - \pi_{k-1} < \delta$ .
12:   Update behaviour policy:  $\theta_{(best)}^k = \arg \max_{j=1, \dots, k} \hat{F}^{\alpha}(j; \pi_j)$ 
13: end for

14: function EVALUATE RISK( $k; \epsilon_k; \beta; P^k; R^k$ )
15:   Set continue-sampling  $\epsilon = 1$ , and  $V = ()$ ;  $r = V = ()$ .
16:   while continue-sampling do
17:     for  $j = 1; 2; \dots$  in parallel do
18:       Sample  $P_j \sim P^k(P|D)$ ,  $R_j \sim R^k(R|D)$ 
19:        $J(\cdot; P_j)$  Policy Evaluation on model  $P_j$ , reward function  $R_j$ 
20:        $V$  append  $J(\cdot; P_j)$ 
21:        $r = V$  append  $r = J(\cdot; P_j)$ 
22:     end for
23:     Sort  $V; r = V_k$  in increasing order of  $V$ . Denote  $N = |V|$ .
24:     Find  $r; s$  with  $j_s = r_j$  minimal, subject to  $\sum_{i=r}^{s-1} \frac{1}{N} \sum_{i=1}^N (1 - \epsilon)^{N-i} > 1$ 
25:     continue-sampling  $\epsilon = V[s] - V[r] \leq \epsilon_{rel}(V[L] - V[1])$ 
26:   end while
27:   Estimate of VaR:  $\hat{F}^{\alpha}(k; \pi_k) = V[bN : c]$ 
28:   CVaR :  $\hat{F}^{\alpha}(k; \pi_k) = \text{mean}(V[1]; \dots; V[bN : c])$ .
29:   return  $V; r = V; \pi_{k+1}; \hat{F}^{\alpha}(k; \pi_k)$ 
30: end function

```

We consider several possibilities for `OptimisticPlanner` in this work but this is not meant to be an exhaustive list nor the best possible that could be achieved. The objectives we consider here will be referred to as “PSRL”, “Upper CVaR”, and “MaxOpt”, which we describe below.

PSRL: The `OptimisticPlanner` ($\cdot; J$) in PSRL (Osband and Van Roy, 2017) is a randomized algorithm that returns $\pi_{P_i} = \arg \max_{\pi} J(\cdot; \pi)$, where $P_i \sim (jD)$. We check whether the constraint is satisfied for π . If not, we resample $P_i \sim (jD)$, and repeat the process. This procedure is computationally very expensive. Every verification of whether the constraint is satisfied for a given policy π_i requires the computation of $\text{CVaR}(J(\cdot; \pi_i))$, which itself requires many samples from the posterior. In addition, this procedure might need to be repeated many times if the constraint is a difficult one to satisfy.

Upper CVaR: We consider the opposite end of the value function distribution from the constraint. We define $\text{CVaR}(Z) = \min_{f} \int_{\alpha}^1 E[(Z - f)^+];$ where $\alpha \in (0; 1)$ is the quantile order defined on the upper tail of the distribution of Z . The `OptimisticPlanner` is then $\text{OptimisticPlanner}(\cdot; J) = \arg \max_{\pi} \text{CVaR}(J(\cdot; \pi))$: Referring to Figure 1, this objective optimizes the average performance in the upper α -quantile of the posterior. Therefore, this objective acts optimistically with respect to the agent’s belief about the true MDP, while the constraint ensures that the agent does not have many bad events given its belief.

Algorithm 2 PG Algorithm for Optimistic Risk-Constrained MBRL

```

1: Input: prior distributions  $P_0(P)$ ;  $R_0(R)$ , quantile order  $\alpha \in (0; 1)$ , relative error tolerance  $\epsilon_{rel}$ , initial state
   distribution  $\mu_0$ , confidence level  $\beta$ , policy learning rates  $\eta$ , convergence criterion  $\delta$ , initial policy parameters  $\theta_0$ .
2: Initialization: policy parameters  $\theta_0^{(0)} = \theta_0$ , Lagrangian parameter  $\lambda = \lambda_0$ , initial penalty parameter  $\rho_0$ .
3: for  $i = 0; 1; 2; \dots$  do
4:   Collect transition(s) from environment and add to data  $\mathcal{D} = \{f, X_t; A_t; X_{t+1}; R_t\}$ 
5:   Update posteriors  $P^{(i)}(P|D)$ ;  $R^{(i)}(R|D)$ 
6:   Set risk constraint  $b = \max_{P \sim P^{(i)}} \text{CVaR}(\mathcal{J}(P))$  (using Alg 1). Re-initialize  $\mu_0, \theta_0$ 
7:   for  $k = 0; 1; 2; \dots$  do
8:     repeat
9:        $V_k; r, V_k; \text{CVaR} = \text{EVALUATE RISK}(\theta_k^{(i)}; \epsilon_{rel}; \beta; P^{(i)}; R^{(i)})$ 
10:      Update:  $\theta_{k+1}^{(i)} = \theta_k^{(i)} + \eta (k) (r - f(\theta_k^{(i)}; P)) + \frac{1}{N_k} \sum_{i=0}^{N_k} (V_k[i] - r) \theta_k^{(i)}$ 
11:      until  $\|\theta_{k+1}^{(i)} - \theta_k^{(i)}\| \leq \delta$ 
12:      Update:  $\lambda_{k+1} = \lambda_k + \frac{1}{k} (\text{CVaR}(\mathcal{J}(\theta_{k+1}^{(i)}; P)) - b)$ 
13:      Penalty Update:  $\rho_{k+1} = \rho_k \cdot 2^{(0; \lambda_k)}$ 
14:    end for
15:    Set policy parameters  $\theta_0^{(i+1)} = \theta_{k+1}^{(i)}$ .
16:  end for

```

MaxOpt: Similarly, in “MaxOpt”, the idea is to act as optimistically as possible with respect to samples from the posterior (maximum sample in Figure 1), instead of taking the UpperCVaR approach where we are “cautiously optimistic”. The objective is then $\text{OptimisticPlanner}(\mu; \mathcal{J}) = \arg \max_{P \in \mathcal{P}} \mathcal{J}(\mu; P)$; $\mathcal{P}_1; \dots; \mathcal{P}_m = \mathcal{P}(\mathcal{J}|D)$; where the number of models m is determined with the stopping condition used to estimate the constraint CVaR. This is simply a practical heuristic as the samples for estimating the constraint can be reused.

If $\text{OptimisticPlanner}(\mu; \mathcal{J})$ can be written as the optimization of a function, we denote $\max_{P \in \mathcal{P}} f(\mu; P)$. To solve (8), we use an Augmented Lagrangian approach (Wright et al., 1999), with multiplier λ and Lagrangian:

$$L(\mu; \lambda) = f(\mu; P) + \lambda (\text{CVaR}(\mathcal{J}(\mu; P)) - b); \quad (9)$$

We present the full algorithm in Alg. 2. At each iteration of the outer loop, we collect data and update posteriors. Given penalty $\rho_k > 0$, at each iteration of the inner loop, we first find an approximate maximizer θ_{k+1} of (9) by taking gradient ascent steps using the gradient of the Lagrangian: $\nabla_{\theta} L(\theta; \lambda) = r - f(\theta; P) + \lambda \text{CVaR}(\mathcal{J}(\theta; P))$; where the first term is a subgradient if θ is “Upper CVaR” or MaxOpt. We then update the multiplier λ to obtain λ_{k+1} using:

$$\lambda_{k+1} = \text{proj}_{\mathbb{R}^+} \left(\lambda_k + \rho_k (r - L(\theta_{k+1}; \lambda_k)) \right) = \text{proj}_{\mathbb{R}^+} \left(\lambda_k + \frac{1}{k} (\text{CVaR}(\mathcal{J}(\theta_{k+1}; P)) - b) \right); \quad (10)$$

where $\text{proj}_{\mathbb{R}^+}$ is a projection operator on \mathbb{R}^+ . Finally, we make the penalty parameter ρ_{k+1} smaller than ρ_k . The starting point for the next iteration is set to the solution of the current iteration. As the penalty parameter becomes smaller, constraint violations are penalized more and the infeasibility in the solution becomes smaller (Wright et al., 1999). Another element of Alg. 2 is in the choice of acceptable risk. Ideally we would like this value to be adaptively chosen, and make sense based on the environment and collected data. As a result, we propose to choose $b = \max_{P \in \mathcal{P}} \text{CVaR}(\mathcal{J}(\mu_0; P))$; with $\mu_0^{(i)}$ the solution from the previous iteration of data collection.

5. Empirical Studies

We present preliminary results on Alg.’s 1 and 2 on some finite domains: Chain (Strens, 2000), DoubleLoop (Strens, 2000), Frozen Lake (FL) (Sutton and Barto, 2018), Cliff-Walking (Sutton and Barto, 2018), LavaLake (Leike et al., 2017), and IslandNavigation (Leike et al., 2017). All these environments require a balance of exploration and exploitation, and other than Chain and DoubleLoop, have some states that are low value and should be avoided. The full

experimental setup we use and descriptions of domains we consider are given in Appendix A. As we are in a discrete state space, we use as prior a Dirichlet distribution with equal parameters $\beta = 0.5$ for all states in X and estimate the likelihood of the data using a categorical distribution, which results in a Dirichlet posterior with parameters $\beta_0 + n_{x;a;x}$, where $n_{x;a;x}$ are the transition counts for each transition from x to a in the collected data. We use a Normal distribution with mean μ and variance σ^2 for the prior of the reward function.

With Alg. 1, we directly optimize CVaR of the policy performance according to the posterior as we collect more data and make updates. Alternatively, we could optimize the expected value with respect to the posterior. We denote this as “PG” for policy gradient, and “PG with CE” for policy gradient with Certainty Equivalent model. In the former, we take the average gradient step over samples of MDPs. In the latter, we take gradient steps with respect to the certainty equivalent, or MLE, model. We expect that the overall performance of optimizing the CVaR would be worse for environments where error in the model does not cause catastrophic drops in performance. We expect optimizing the risk would perform worse as it is more cautious with respect to different realizations of the true environment, potentially exploring the environment less.

The results for comparing the CVaR, PG, and PG with CE are shown in part (b) of Figure 2 (and 4, 5, 6, 7, 8 in Appendix). We can see that across all environments tested, the PG with CE model has the worst CVaR on the posterior, indicating that although it may have good performance in certain cases, it is not acting safely with respect to its uncertainty in the environment, as expected. In addition, as expected, for most environments the PG or PG with CE perform better than CVaR. However, in environments that require some degree of “safe” or conservative behaviour (i.e. LavaLake, IslandNavigation) CVaR outperforms the other methods.

We also evaluate Alg. 2. As the constraint here limits the space of policies, simply planning optimistically without a constraint may result in better overall performance on the true environment as more optimistic policies may violate the constraint during training. On the other hand, as before, if the environment has areas of low-probability but very low-value, in the early stages of training our agent may not have collected much or any data from these areas, and thus it will contain inaccuracies (due to the prior) and uncertainty about these areas. In this case, having the risk constraint over our belief should lower the chances of learning a policy during planning that would venture into these areas of the true environment as the policy is in general more cautious with respect to its belief.

First, we check the effect of adding a risk constraint to our existing PG planner. The results of this experiment are shown in part (d) of Figure 2 (and 4, 5, 6, 7, 8 in Appendix). In parts (a) and (c) of the same Figures, we show the results of adding a risk constraint to an optimistic planner.

For the PG variants, again we notice a marked improvement in the policy performance in more “risky” environments (FrozenLake, LavaLake, and IslandNavigation) for the constrained version, and similar performance to the unconstrained versions in the other environments. For the optimistic planner variants (parts (a) and (c) of the Figures), we find a similar pattern as the PG variants (part (d)), although we also see the unconstrained planners sometimes outperforming the constrained versions in both performance and CVaR. This could point to an issue with optimization, as optimization of the constrained problem requires us to first find the optimal CVaR (Alg. 2) and also to find constraint-feasible policies, which may be more difficult if the constraint is active.

Another observation is that in general, as the methods learn better policies, their CVaR also increases. This points to the idea that a good, high-performing policy naturally will not have low CVaR as catastrophic events are already minor. However, we also do not have any guarantees that the policy behaves safely with respect to our belief about the environment, including during learning. Being risk-aware may eliminate intermittent policies that may in expectation have good performance but also have low CVaR.

6. Related Works

Taking parameteric uncertainty into account in model-based methods has been extensively studied in the literature. For a full review of the literature in this field, up to the year 2015, we refer the reader to Chapter 6 of (Ghavamzadeh et al., 2015). We highlight the most relevant as well as more recent works. Delage and Mannor (2010) investigate the “percentile” criterion, which, in the context of this work, is the VaR (applied to the posterior distribution). For the case of a Dirichlet prior, Delage and Mannor (2010) show that an approximately optimal solution can be derived. The setting of Delage and Mannor (2010) is for a fixed dataset and the exploration of a policy as more data is collected is not considered. Additionally, our approach is based on PG.

Figure 2: IslandNavigation. Solid lines are Avg of 8 runs, shaded areas Std. Err. (Top) Env Returns, (Bottom) Posterior CVaR.

Chen and Bowling (2012) generalize the percentile criterion by replacing it with a measure over percentiles. For example, the CVaR can be obtained by using a specific measure over percentiles. Chen and Bowling (2012) also propose a family of percentile measures “k-of-N” that can be tuned to closely approximate the CVaR. This measure is equivalent to sampling in the following way: for a policy sample π , select the set of k MDPs with the worst performances $\{P_i\}_{i=1:k}$, and choose an MDP uniformly at random from this set. The authors show that the k-of-N measure can be solved efficiently and give convergence guarantees. We point out two key differences between Alg. 1 and the approach of Chen and Bowling (2012). We use PG instead of counterfactual regret minimization. This allows us to easily incorporate the risk-aware objective as a constraint (Section 4). The second difference is in the sampling procedure used to approximate, for example, the CVaR, as k-of-N could also be used in PG calculations. Comparing the way we estimate the CVaR (in Alg. 1) to the k-of-N sampling procedure, we note that for comparable k and N parameters, our approach has a factor of k lower variance in estimating the PG than k-of-N since the latter uses a single sample to estimate the PG whereas in our approach, the PGs for around k samples would be averaged.

The problems we consider in this work are sometimes addressed under the formulation of Bayes-Adaptive MDP (BAMDP) (Duff, 2002), where the problem is cast into a belief-state MDP by augmenting the state with the posterior belief. The optimal policy in this case (i.e. the Bayes-optimal policy) makes the optimal tradeoff between exploration and exploitation with respect to its belief (Duff, 2002). In our work, we aim to explicitly prevent the agent from taking actions in likely-to-be unfavourable areas. In other words, we sacrifice Bayes-optimality in order to ensure the policies during learning are “safe” with respect to the parametric uncertainty. Sharma et al. (2019) consider the problem of trading off exploration and risk-sensitivity over model space in Bayesian RL. The formulation in this work is different from ours as they formulate the problem as a two-player zero-sum game between the agent and an adversary that tries to modify the agent’s beliefs. They also use a different approach to solving the problem instead of PG.

7. Conclusions and Future work

We introduced two algorithms based on applying the notion of risk to address parametric uncertainty in model-based RL. One was based on directly optimizing the CVaR using PG. The other was using the CVaR as a constraint, and acting optimistically in the environment to collect data and explore. We empirically evaluated these algorithms on some finite domains and saw promising results. As part of future work, we would like to extend this work and the empirical evaluations to continuous domains with neural network function approximators, and posterior approximations. PG methods have been applied successfully in these domains (Sutton et al., 1999; Schulman et al., 2015, 2017) and thus this is a natural next step for our approach. Note that since we do not exactly calculate the CVaR, but estimate them using a finite number of samples and with a given confidence level, there is an error introduced in the calculation of the subgradient. It is not known how this error affects the gradient step procedure and its convergence, and we leave analysis of this to future work. For the constrained algorithm 2, an important question is how to handle situations where the constraint is difficult to achieve. We also leave exploration of further optimistic objectives for future work.

8. Acknowledgements

AMF acknowledges the funding from the Canada CIFAR AI Chairs program, as well as the support of the Natural Sciences and Engineering Research Council of Canada (NSERC) through the Discovery Grant program. Resources used in preparing this research were provided, in part, by the Province of Ontario, the Government of Canada through CIFAR, and companies sponsoring the Vector Institute.

References

- Carlo Acerbi and Dirk Tasche. Expected shortfall: a natural coherent alternative to value at risk. *Journal of Economic Surveys* 31(2):379–388, 2002.
- Giorgio Angelotti, Nicolas Drougare, and Caroline Ponzoni Carvalho Chanel. Exploitation vs caution: Risk-sensitive policies for of ine learning. *arXiv preprint arXiv:2105.13431*, 2021.
- András Antos, Csaba Szepesvári, and Rémi Munos. Learning near-optimal policies with Bellman-residual minimization based tted policy iteration and a single sample path. *Machine Learning* 71:89–129, 2008.
- Philippe Artzner, Freddy Delbaen, Jean-Marc Eber, and David Heath. Coherent measures of risk. *Mathematical Finance*, 9(3):203–228, 1999.
- Karl Johan Astöm and Björn Wittenmark. *Adaptive Control*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2nd edition, 1994.
- Dimitri P. Bertsekas and John N. Tsitsiklis. *Neuro-Dynamic Programming*. Athena Scientific, 1996.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Nectoula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL: <http://github.com/google/jax>.
- K Briggs and FM Ying. How to estimate quantiles easily and reliably. *Mathematics Today* 2018(February), 2018.
- Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- Jinglin Chen and Nan Jiang. Information-theoretic considerations in batch reinforcement learning. *Proceedings of the 36th International Conference on Machine Learning (ICML)* 2019.
- Katherine Chen and Michael Bowling. Tractable objectives for robust policy optimization. *Advances in Neural Information Processing Systems* 25, 2012.
- Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research* 18(1):6070–6120, 2017.
- Erick Delage and Shie Mannor. Percentile optimization for markov decision processes with parameter uncertainty. *Operations research* 58(1):203–213, 2010.
- Pierluca D’Oro, Alberto Maria Metelli, Andrea Tirinzoni, Matteo Papini, and Marcello Restelli. Gradient-aware model-based policy search. *Thirty-Fourth AAAI Conference on Artificial Intelligence (AAAI)* 2020.
- Michael O’Gordon Duff. *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes*. University of Massachusetts Amherst, 2002.
- Benjamin Eysenbach, Alexander Khazatsky, Sergey Levine, and Ruslan Salakhutdinov. Mismatched no more: Joint model-policy optimization for model-based RL. *arXiv preprint arXiv:2110.02758*, 2021.
- Amir-massoud Farahmand, Mohammad Ghavamzadeh, Csaba Szepesvári, and Shie Mannor. Regularized policy iteration with nonparametric function spaces. *Journal of Machine Learning Research (JMLR)* 17(139):1–66, 2016.

- Amir-massoud Farahmand, Andre Barreto, and Daniel Nikovski. Value-aware loss function for model-based reinforcement learning. In *Artificial Intelligence and Statistics* pages 1486–1494. PMLR, 2017.
- Mohammad Ghavamzadeh, Shie Mannor, Joelle Pineau, Aviv Tamar, et al. Bayesian reinforcement learning: A survey. *Foundations and Trends® in Machine Learning* 6(5-6):359–483, 2015.
- Christopher Grimm, Andre Barreto, Satinder Singh, and David Silver. The value equivalence principle for model-based reinforcement learning. *Advances in Neural Information Processing Systems* 33:5541–5552, 2020.
- László Györfi, Michael Kohler, Adam Krzyżak, and Harro Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer Verlag, New York, 2002.
- David Ha and Jürgen Schmidhuber. World models. *arXiv preprint arXiv:1803.10122*, 2018.
- Danijar Hafner, Timothy Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. *arXiv preprint arXiv:1912.01603*, 2019.
- Garud N. Iyengar. Robust dynamic programming. *Mathematics of Operations Research* 30(2):257–280, 2005.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research (JMLR)* 11:1563 – 1600, 2010.
- Nathan Lambert, Brandon Amos, Omry Yadan, and Roberto Calandra. Objective mismatch in model-based reinforcement learning. *arXiv preprint arXiv:2002.04523*, 2020.
- Jan Leike, Miljan Martic, Victoria Krakovna, Pedro A Ortega, Tom Everitt, Andrew Lefrancq, Laurent Orseau, and Shane Legg. Ai safety gridworlds. *arXiv preprint arXiv:1711.09883*, 2017.
- Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research (JMLR)* 9:815–857, 2008.
- Arnav Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain transition matrices. *Operations Research* 53(5):780–798, 2005.
- Ian Osband and Benjamin Van Roy. Why is posterior sampling better than optimism for reinforcement learning? In *International conference on machine learning*, pages 2701–2710. PMLR, 2017.
- Ian Osband, Daniel Russo, and Benjamin Van Roy. (more) efficient reinforcement learning via posterior sampling. In *Advances in Neural Information Processing Systems*, volume 26, 2013.
- R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of Risk* 2:21–42, 2000.
- R Tyrrell Rockafellar, Stan Uryasev, and Michael Zabarankin. Master funds in portfolio analysis with general deviation measures. *Journal of Banking & Finance* 30(2):743–778, 2006.
- Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, et al. Mastering atari, go, chess and shogi by planning with a learned model. *Nature* 588(7839):604–609, 2020.
- John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897. PMLR, 2015.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, 2009.

- Apoorva Sharma, James Harrison, Matthew Tsao, and Marco Pavone. Robust and adaptive planning under model uncertainty. In Proceedings of the International Conference on Automated Planning and Scheduling, volume 29, pages 410–418, 2019.
- Ingo Steinwart and Andreas Christmann. Support Vector Machines. Springer, 2008.
- Malcolm Strens. A bayesian framework for reinforcement learning. In International Conference on Machine Learning (ICML), pages 943–950, 2000.
- Richard S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In Proceedings of the 7th International Conference on Machine Learning (ICML), 1990.
- Richard S Sutton and Andrew G Barto. Reinforcement learning: An introduction. MIT press, 2018.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. Advances in neural information processing systems, 1999.
- Csaba Szepesvári. Algorithms for Reinforcement Learning. Morgan Claypool Publishers, 2010.
- Aviv Tamar, Yonatan Glassner, and Shie Mannor. Optimizing the CVaR via sampling. Twenty-Ninth AAAI Conference on Artificial Intelligence, 2015.
- Samuele Tosatto, Matteo Pirotta, Carlo D'Eramo, and Marcello Restelli. Boosted fitted Q-iteration. Proceedings of the 34th International Conference on Machine Learning (ICML), 2017.
- Stephen Wright, Jorge Nocedal, et al. Numerical optimization. Springer Science & Business Media, 1999.
- Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. preprint arXiv:1910.08348, 2019.

Appendix A. Experimental Setup

A.1 Environment Descriptions

A.1.1 CHAIN

This is the environment defined by [Strens \(2000\)](#) and shown in Fig. 3b. There are two actions and 5 states. There is also a 0.2 probability of the agent's intended action being changed to the other action. Once a state is reached, the agent stays there and receives a reward until it "slips" again with probability 0.2. The optimal policy is to always choose action a , even though sometimes this results in the other action being taken. This problem requires effective exploration to solve.

A.1.2 DOUBLELOOP

This environment is also defined by [Strens \(2000\)](#) and shown in Fig. 3c. It involves two loops of length 5 joined at a single start state. Two actions are available and transitions are deterministic. Taking action a repeatedly causes traversal of the right loop, with a reward of 1 for every 5 actions taken. Taking action b repeatedly causes traversal of the left loop, with a reward of 2 for every 5 actions taken. This problem requires a difficult compromise between exploration and exploitation.

A.1.3 CLIFFWALKING

This is the environment proposed by [Sutton and Barto, 2018](#)), and illustrated in Fig. 3a. The Goal state is an absorbing state with reward 1. Each action has a 0.2 probability of "failing", and taking the agent to a random adjacent state.

A.1.4 FROZEN LAKE (FL)

This is the Frozen Lake environment of OpenAI Gym [Brockman et al. \(2016\)](#) with modified rewards. In this environment, the agent moves in a 4x4 grid starting in the upper left corner towards an absorbing goal state that has a reward of 10 and is in the lower right corner of the grid. Along the way, the grid has some holes, which are also absorbing states with rewards of 1. The rest of the grid has rewards sampled uniformly from the interval $(0, 0.8)$ at the generation of the map, and the rewards are fixed thereafter (deterministic). The grid is also "slippery" with probability $\frac{1}{3}$ the agent will move in the intended direction with probability $\frac{1}{3}$ and will move in either perpendicular direction with equal probability of $\frac{1}{3}$ in both directions. The map is randomly generated at the start of each episode, ensuring there is a hole-free path from the start to the goal.

A.1.5 LAVA LAKE AND ISLAND NAVIGATION

These domains are taken from [Leike et al. \(2017\)](#) and shown in Figures 3d and 3e. They each have 4 actions and are episodic, resetting to a randomized starting position (the figures). The agent must reach a Goal state upon which it receives a reward of 50 and the episode ends. Otherwise, a reward of 0 is incurred unless the agent lands in the Lava/Water, with a reward of -50 and termination of the episode. In the IslandNavigation domain additional information is provided to the agent through a "safety distance" that maps the current environment state to the agent's Manhattan distance to the closest water cell ([Leike et al., 2017](#)).

A.2 Experimental Details

The policy is represented using parameters $(\theta_0; \dots; \theta_{x_a}; \dots; \theta_{j \times j \times j})$ for $x_a \in (X \setminus A)$, and defined as:

$$(\pi_j)_a = \frac{P_{\theta}^{x_a}}{\sum_{a'} P_{\theta}^{x_{a'}}}$$

(a) Cliff-walking reproduced from [Sutton and Barto \(2018\)](#).

(b) Chain MDP reproduced from [Strens \(2000\)](#).

(c) Chain MDP reproduced from [Strens \(2000\)](#).

(d) LavaLake MDP reproduced from [Leike et al. \(2017\)](#)

(e) IslandNavigation MDP reproduced from [Leike et al. \(2017\)](#)

Figure 3: Illustrations for two of the MDPs used in the experiments.

Before policy updates begin, the posteriors are updated with data collected from a random policy for 10,000 environment steps.

The PG is calculated exactly on collected samples from the posterior, as we are in a tabular setting (i.e. we do not use samples from each model separately but rather solve each MDP exactly).

Hyperparameters, for all envs:

- quantile order : 0:1,
- upper quantile order : 0:1,
- significance level : 0:1,
- relative error tolerance ϵ_{rel} : 0:1,
- discount factor: 0.99,
- policy learning rate: 0.1 for unconstrained algorithms, 0.01 for constrained,
- learning rate: 0.01,
- Number of iterations of constrained optimization (outer loop, i.e. number of updates and penalty): 20,
- Number of samples taken for PG: 500, for risk-based methods the number could be adaptively chosen based on confidence intervals, but we found a fixed number 500 was both sufficient for the given confidence interval and efficient.
- Seeds are set in the numpy random number generator, the jax [Bradbury et al. \(2018\)](#) random number generator, and to seed the environments. It would be possible that two runs with the same seeds are not exactly identical however since there may be lower level sources of randomness we did not account for.

Appendix B. Additional Experiments

We present additional experiments in Figures [4](#), [5](#), [6](#), [7](#), [8](#).

(a) Comparison of Max Opt algorithm and Max Opt with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

(b) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and PG with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

(c) Comparison of Upper CVaR algorithm and Upper CVaR with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

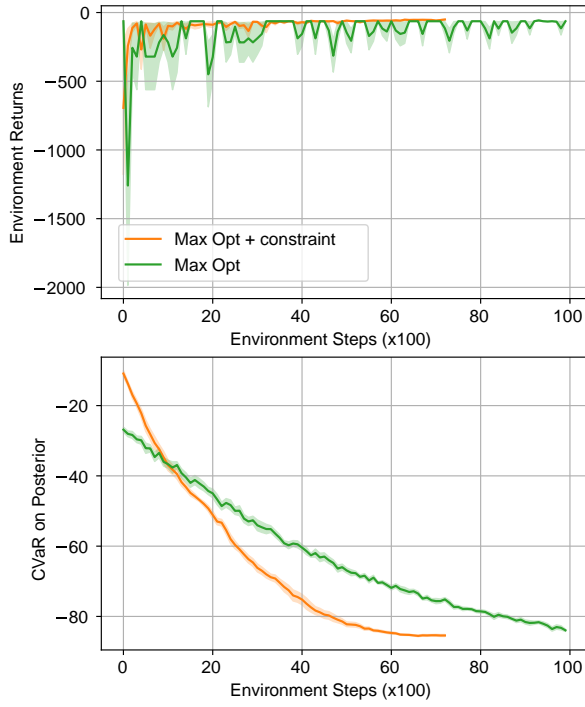
(d) Comparison of PG algorithm and PG with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

(a) Comparison of Max Opt algorithm and Max Opt with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

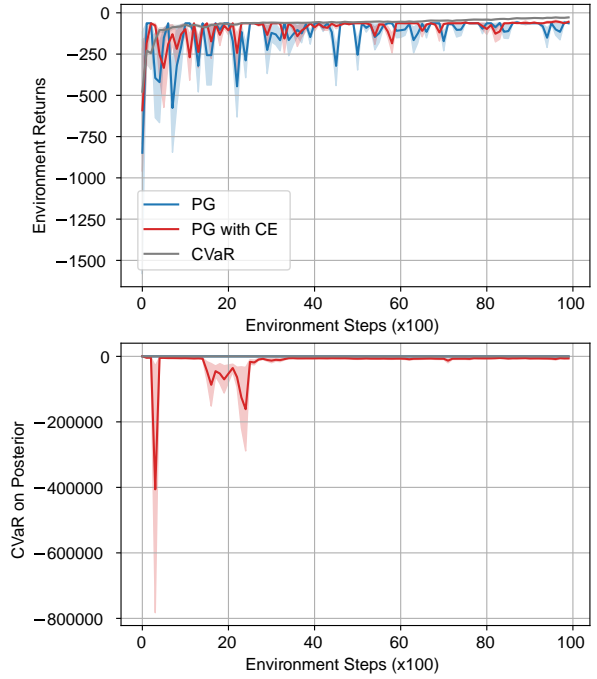
(b) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and PG with the Certainty Equivalent or MLE model, and optimizing the CVaR directly. (Top) Average Rewards, (Bottom) CVaR on Posterior.

(c) Comparison of Upper CVaR algorithm and Upper CVaR with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

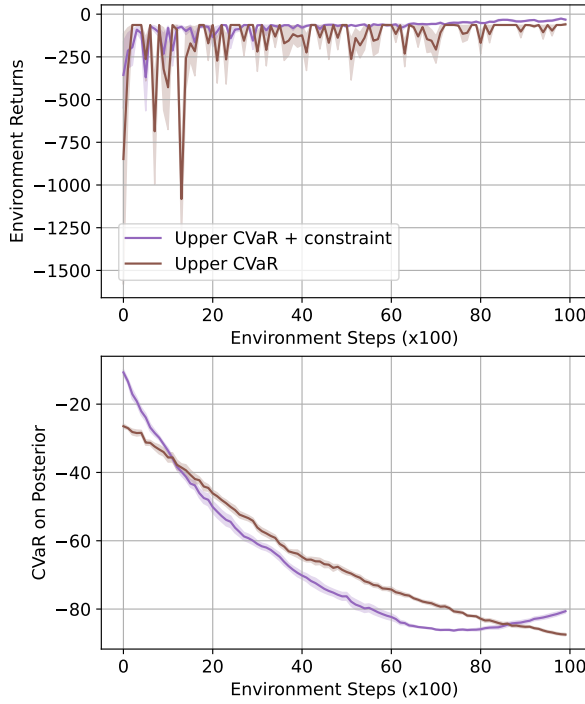
(d) Comparison of PG algorithm and PG with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.



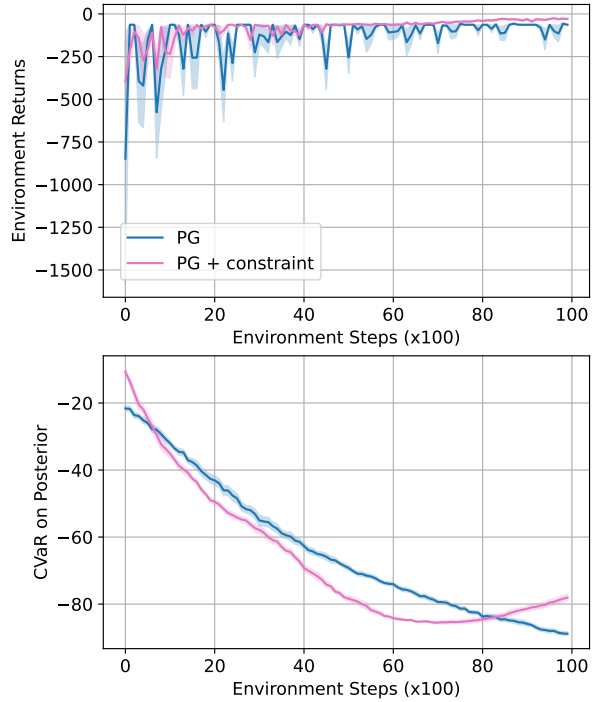
(a) Comparison of Max Opt algorithm and Max Opt with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.



(b) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and optimizing the CVaR directly. (Top) Average Rewards, (Bottom) CVaR on Posterior.

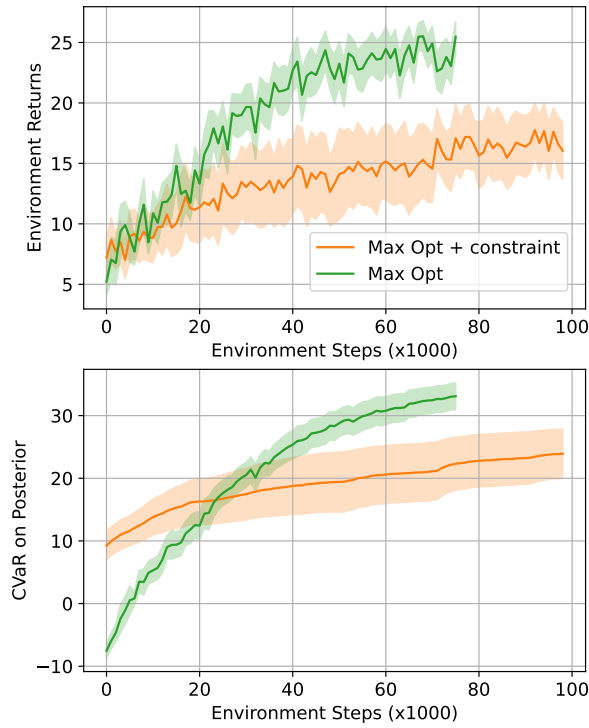


(c) Comparison of Upper CVaR algorithm and Upper CVaR with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

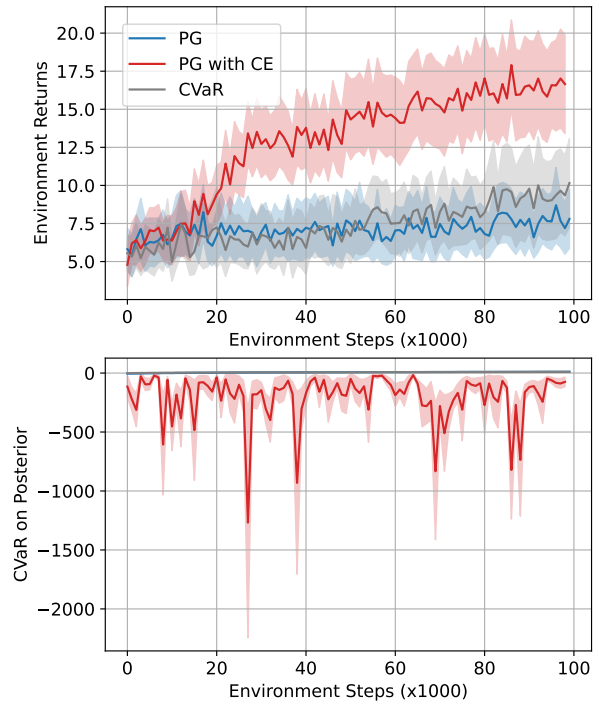


(d) Comparison of PG algorithm and PG with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

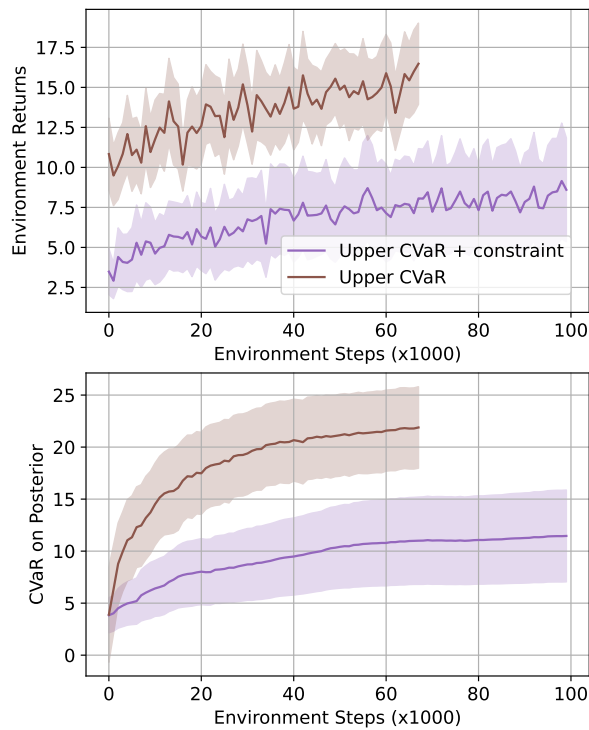
Figure 6: Performances on CliffWalking env



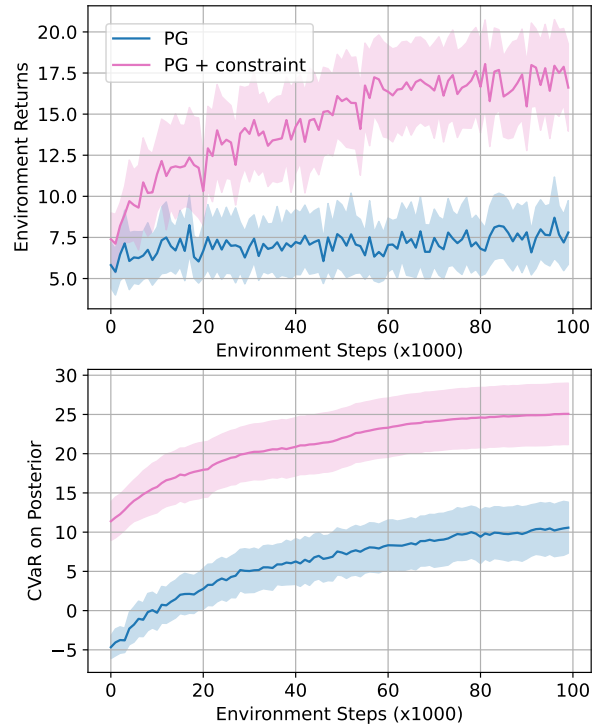
(a) Comparison of Max Opt algorithm and Max Opt with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.



(b) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and optimizing the CVaR directly. (Top) Average Rewards, (Bottom) CVaR on Posterior.

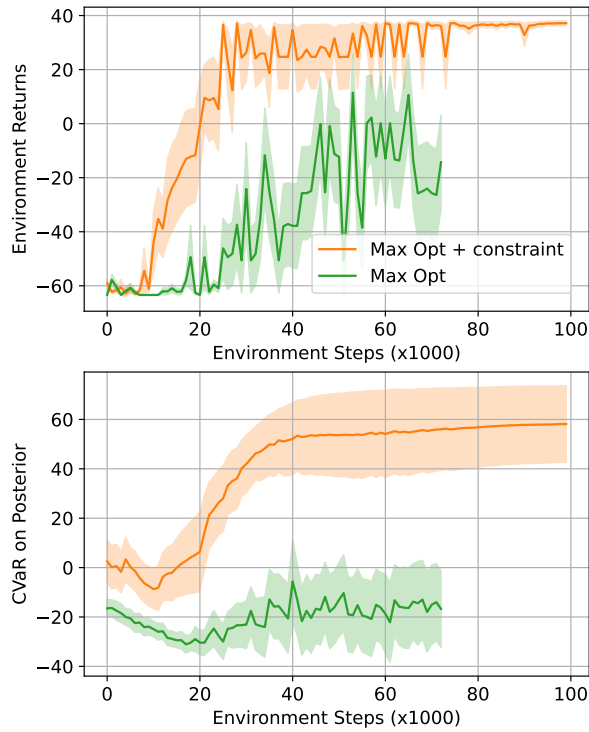


(c) Comparison of Upper CVaR algorithm and Upper CVaR with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.

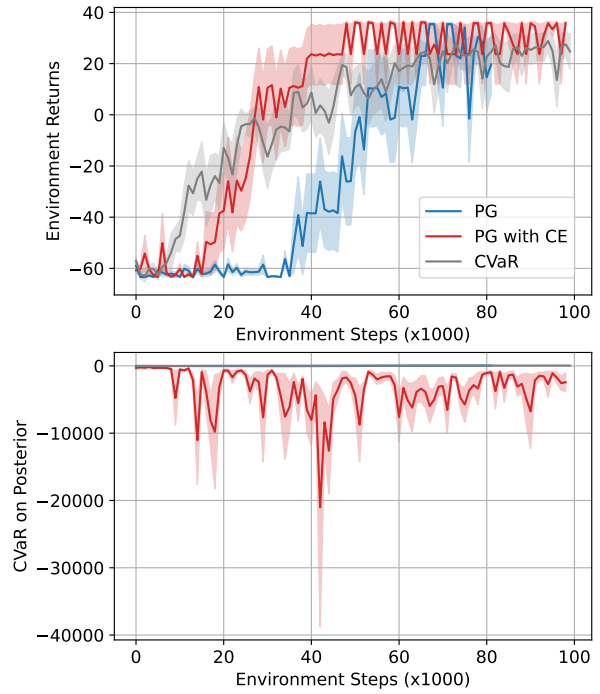


(d) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and optimizing the CVaR directly. (Top) Average Rewards, (Bottom) CVaR on Posterior.

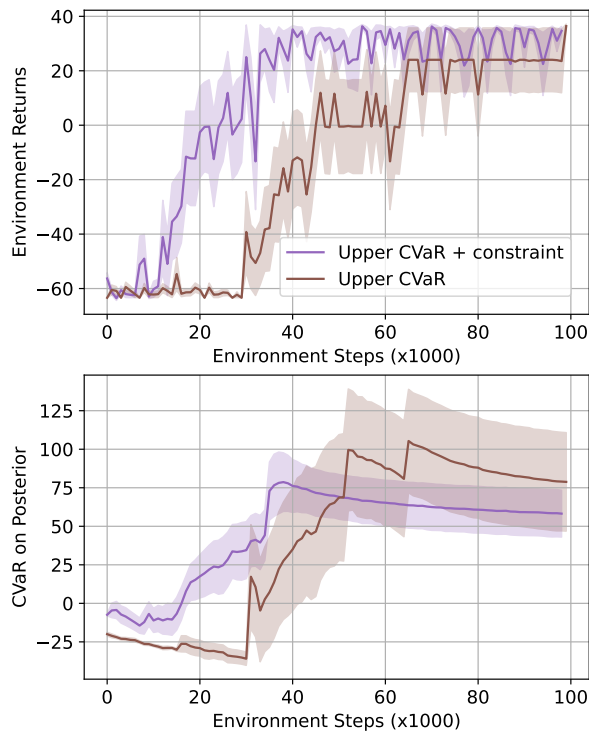
Figure 7: Performances on FrozenLake (FL) env



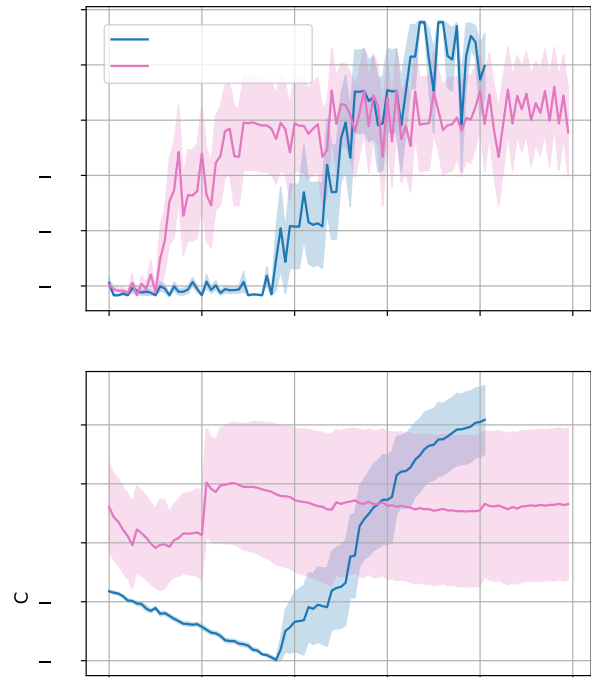
(a) Comparison of Max Opt algorithm and Max Opt with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.



(b) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and optimizing the CVaR directly. (Top) Average Rewards, (Bottom) CVaR on Posterior.



(c) Comparison of Upper CVaR algorithm and Upper CVaR with CVaR-constraint. (Top) Average Rewards, (Bottom) CVaR on Posterior.



(d) Comparison of PG algorithm using samples from the posterior, PG with the Certainty Equivalent or MLE model, and optimizing the CVaR directly. (Top) Average Rewards, (Bottom) CVaR on Posterior.

Figure 8: Performances on LavaLake