# An Offline Risk-aware Policy Selection Method for Bayesian Markov Decision Processes

Giorgio Angelotti[a,b,*], Nicolas Drougard[a,b], Caroline P. C. Chanel[a,b]

[a]*ANITI - Artificial and Natural Intelligence Toulouse Institute, University of Toulouse, France*
[b]*ISAE-SUPAERO, University of Toulouse, France*

**Abstract**

In Offline Model Learning for Planning and in Offline Reinforcement Learning, the limited data set hinders the estimate of the Value function of the relative Markov Decision Process (MDP). Consequently, the performance of the obtained policy in the real world is bounded and possibly risky, especially when the deployment of a wrong policy can lead to catastrophic consequences. For this reason, several pathways are being followed with the scope of reducing the model error (or the distributional shift between the learned model and the true one) and, more broadly, obtaining risk-aware solutions with respect to model uncertainty. But when it comes to the final application which baseline should a practitioner choose? In an offline context where computational time is not an issue and robustness is the priority we propose Exploitation vs Caution (EvC), a paradigm that (1) elegantly incorporates model uncertainty abiding by the Bayesian formalism, and (2) selects the policy that maximizes a risk-aware objective over the Bayesian posterior between a fixed set of candidate policies provided, for instance, by the current baselines. We validate EvC with state-of-the-art approaches in different discrete, yet simple, environments offering a fair variety of MDP classes. In the tested scenarios EvC manages to select robust policies and hence stands out as a useful tool for practitioners that aim to apply offline planning and reinforcement learning solvers in the real world.

*Keywords:* Risk-aware Markov Decision Process, Bayesian Markov Decision Process, Offline Policy Selection, Offline Model Learning for Planning, Offline Reinforcement Learning

## 1. Introduction

The deployment of autonomous agents in an unknown, and possibly stochastic, environment is a delicate task that usually requires a continuous agent-environment interaction which is not always affordable in real-life situations. For instance, in applications such as the training of medical robots and automated vehicles [1, 2] the interaction with the environment can be both too risky (e.g. proximity to a human) and costly since: (i) any mistake could lead to catastrophic aftermaths; or (ii) the data collection phase requires a direct human involvement, which is usually expensive and time demanding. Hence, it can be convenient to exploit previously collected data sets in order to limit additional (dangerous or superfluous) interaction.

---

*Corresponding author

Offline model learning for planning and offline Reinforcement Learning (RL) are the branches of machine learning that leverage previously collected batches of experiences with the aim of establishing an optimal behavioral policy offline. In recent years, the RL community published a great number of papers on the subject, as for instance, the works in [3, 4, 5, 6, 7, 8, 9, 10, 11], demonstrating the growing interest in the field. The proposed algorithms try to improve the performance of a policy obtained either with model-free or model-based RL approaches. The intuition behind these methods is always the same: optimizing a trade-off between *exploitation and caution*. The policy optimization procedure is usually tailored in order to generate strategies that are not too distant from the one originally used to collect the batch. In this way, the agent will follow a strategy that will not drive him towards regions of the state-action space for which it possesses a high degree of uncertainty.

For instance, in the offline RL literature tailored policy optimization procedures for Markov Decision Processes (MDPs) are implemented: (i) in Conservative Q-Learning (CQL) approach [12], by limiting the overestimation of Out-Of-Distribution transitions; (ii) in the Behavior Regularized Actor-Critic (BRAC) paradigm [6], as a penalty in the value function proportional to an estimate of the policies' distributional shift; (iii) in the Model-based Offline Policy Optimization (MOPO) algorithm [9], as a penalization added to the reward function which is proportional to an estimate of the distributional shift in the dynamical evolution of the system - also called *epistemic* (model) error; and, (iv) in the Model-Based Offline Reinforcement Learning (MOReL) approach [10], by creating an additional and highly penalized absorbing state and by forcing the agent to transit to it when the model uncertainty for a specific state-action pair is above a given threshold.

On top of a non-trivial estimate of per-transition uncertainty, which is often performed with Deep Neural Networks, the said baselines notably require the fine-tuning of domain-dependent hyperparameters [13]. Such an empirical calibration demands additional interaction with the environment and thus betrays the original purpose of offline learning.

How to select the best set of hyperparameters for offline RL baselines based on Deep Neural Networks? This question does not have a trivial answer [11]. Usually, the policies obtained are evaluated offline with *off policy-evaluation* (OPE) approaches like Fitted Value Iteration [14] or Fitted Q Evaluation (FQE) [15]. Indeed, [13] showed that ORL baselines using DNN are not robust with respect to hyperparameter selection and presented a comparison between offline hyperparameter selection methods. Nevertheless, algorithms like FQE require hyperparameters to be tuned, shifting then the problem of selecting the best hyperparameters for the ORL baseline to the one of selecting the hyperparameters for the OPE. Recently, a totally hyperparameter-free method called BVFT for OPE has been proposed in the work in [16]. However, the said method is affected by limited data efficiency and computational complexity that scales quadratically with the number of models to compare. Concurrently, the work in [17] presented a pessimistic method to estimate and select models for Offline Deep RL.

From another perspective, the Safe Policy Improvement with Baseline Bootstrapping (SPIBB) methodology [3] and the Batch Optimization of Policy and Hyperparameter (BOPAH) approach [7] were developed with the aim of obtaining a robust policy with a true hyperparameter agnostic approach. The former generates a safe policy improvement over the data collecting policy with theoretical guarantees similar to the ones achievable in Probably Approximately Correct approaches, the latter uses a classic batch (offline) RL approach with a gradient-based optimization of the hyperparameters using held-out data.

Aside from the RL community, researchers whose field is mostly offline model learning for planning deal with the problem of solving MDPs under model uncertainty by focusing on the resolution of *robust MDPs* [18, 19]. In this context, the model dynamics (e.g. a stochastic

2

transition function) lie in a constrained ambiguity set which is a subset included in the whole set of distributions. The problem is hence formulated as a dynamic game against a malevolent nature which at every time step chooses the worst model in the set according to the agent action. Subsequently to these works, the research in [20] introduced the *chance constrained MDP* approach which optimizes policies for the *percentile criterion*: the Value-at-Risk (*VaR*) metric. Reference [20] proved that robust MDPs can generate overly conservative strategies depending on the size and shape of the ambiguity set. And, reference [21] proposed approximate solutions to generate safe policy improvements of the data collector policy. Recently, the works in [22, 23] incorporated prior knowledge upon a Bayesian methodology in order to obtain less conservative ambiguity sets that can yield tighter safe returns estimates. In particular, the study in [22] proposes *Bayesian Credible Region* (BCR), an algorithm that constructs ambiguity sets from Bayesian credible (or confidence) regions and uses them to optimize the risk-aware problem (the robust MDP). Reference [24] introduced the *Soft-Robust Value Iteration* (SRVI) algorithm to optimize for the *soft-robust criterion*, a weighted average between the classic value function and the Conditional *VaR* risk metric with epistemic model uncertainty, to solve a Robust MDP.

In parallel to these works, reference [25] showed that planning in an MDP context using a discount factor $\gamma^*$ lower than the one used in the final evaluation phase $\gamma_{ev}$ yields more performing policies when a trivially learned MDP model is considered. A trivially learned MDP model is said to be the one that maximizes the likelihood of the transitions collected in the batch. Nevertheless, the mathematical expression that should be optimized in order to find $\gamma^*$ is intractable. In the work in [25], the optimal discount factor is finally found by cross-validation which requires additional interaction with the true environment.

Off-policy evaluation for finite, discrete MDPs usually resorts to techniques based on Importance Sampling [26, 11]. The Importance Sampling procedure assigns different weights to samples when one exploits them to estimate values from a distribution that is different from the one that was used to generate the samples. This weight, called the Importance Sampling ratio, in off-policy evaluation is the probability of sampling a specific trajectory using the new policy over the probability of obtaining the same trajectory using the policy used to collect the data. This ratio is independent of the models' transition function and can be simplified as the ratio between the probabilities of generating that given sequence of actions while deploying the two different policies. The Universal Off-Policy Evaluation (UnO) [27] has been proposed to estimate not only the average value and the variance of policy performance but also risk-sensitive metrics based on quantiles like the *VaR* or the *CVaR*. UnO is a non-parametric and model-free estimator based on Importance Sampling for the cumulative distribution of returns of a fixed policy $\pi$ starting from a pre-collected batch of experiences. Estimating the full cumulative distribution allows computing risk-sensitive metrics like the Value at Risk and the Conditional Value at Risk.

Unfortunately, UnO and other Importance Sampling based techniques manage to properly estimate the MDP's value function only for stochastic policies while many policies generated by state-of-the-art approaches are deterministic. On one hand, it is true that a deterministic policy is just a specific type of stochastic policy, but on the other hand, the computation of the importance sampling ratio for deterministic policies collapses onto a Kronecker delta. Eventually, the offline evaluation of deterministic policies with Importance Sampling ratio is not accurate. There is a necessity of developing a technique to evaluate and select offline robust deterministic policies.

Concurrently, with the aim of finding a policy that optimizes the trade-off between *exploitation and exploration* in an *online* setting, model uncertainty has been included in a Bayesian extension of the MDP framework called Bayesian (Adaptive) MDP (BAMDP) [28]. Fixing a prior for the distribution over transition models, a posterior distribution is computed from the likelihood of the

3

sampled trajectories. Some years later, the work in [29] suggested that risk-aware utility functions can replace the common BAMDP value function. In doing so, the said work proposed an algorithm that trades off exploration, exploitation, and robustness (caution). The said works (e.g. [28, 29]) deal with an online context when interaction with the environment is always possible while, as we stated before in this manuscript, we will tackle an offline problem. The works in [30, 31] exploit Bayesian Neural Networks with latent variables to encode the uncertainty in model-based Reinforcement Learning with the generator model represented as a Deep Neural Network. On top of this, optimization for a risk-sensitive utility function that is the future expected return plus a variance term that includes both noise and model uncertainty is performed. Unfortunately, the said risk-sensitive objective can unlikely mitigate efficiently the risk in environments where the distribution of returns is far from Gaussian, *e.g.* multi-modal distributions.

Solving a Bayesian MDP is NP-hard [32]. Moreover, the space $\mathbb{M}$ over which the posterior distribution of models is defined is an infinite set. To overcome this computational constraint, the work in [32] proposes to find the policy that maximizes a utility criterion over a finite list of models. The said framework is called Multi-model MDP. While the maximization of the utility criterion using a finite list of models is more treatable than solving a Bayesian MDP, constraining a possibly infinite set of models to a set of candidate ones could in many cases be a too strong approximation. In this work, we will not limit the space of possible models, but instead of solving a Bayesian MDP, we will rather focus on the problem of selecting the most performant policy from a set of candidates, according to a risk-sensitive criterion, over the full distribution of possible models.

Interestingly, the problem of selecting the right algorithm for a machine learning problem was studied in the work in [33] that formalized the protocol to be followed to analyze and solve an algorithm selection problem. And, following the work in [33], Upper Confidence Bound (UCB-1) for Multi-Armed Bandits has been applied in [34] for algorithm (and hence policy) selection in Reinforcement Learning, however, the said approach was not risk-aware and was limited to an online setting.

In fact, in light of the current limitations of state-of-the-art, we notice that selecting a robust policy for an offline data-driven MDP taking into account the uncertainty in the model learning phase is still an open problem. In this context, our work presents Exploitation vs Caution (EvC), an *offline* Bayesian paradigm to evaluate the performance of the policies provided by the state-of-the-art solvers and *select* the best policy between a set of candidates according to different risk-metrics (*VaR* and also the Conditional *VaR*). The set of candidate policies is initially obtained using the available baselines. Taking inspiration from the works in [25] and [29], the set of candidate policies is then enriched considering strategies that are obtained by solving several MDPs with different discount factors and different transition functions which are sampled from the Bayesian posterior inferred from the original fixed batch of pre-collected transitions. The distribution of the performance of every policy taking into account the model uncertainty is evaluated by alternating Monte Carlo model sampling and Policy Evaluation. In the end the best policy according to the risk-sensitive criterion is selected.

4

In summary, the contributions of this work are:

1. it proposes a paradigm to select the best risk-aware policy between a set of candidate policies for a Bayesian MDP in the offline setting;

2. it gives probabilistic guarantees on the performance of the selected policy by computing risk-aware criteria through a confident estimation method for a given risk level (or quantile order);

3. the empirical results demonstrate the validity of this approach in toy environments compared to the current policy selection baseline (UnO).

The paper is organized as follows: Section 2 starts with a recap of the MDP and of the Bayesian MDP (BMDP) formalisms; risk-aware measures following the prescriptions of the research in [35] and the Risk-aware BMDP are introduced in Section 3; in Section 4 the EvC approach is presented; then in Section 5, the policy selected by EvC is compared against selected baseline approaches, and to a risk-sensitive policy selection approach achieved by evaluating the set of candidate policies with UnO; Section 6 concludes the paper by discussing its limitations and pointing to future work perspectives.

## 2. Background

**Definition 1.** *A Markov Decision Process (MDP) is a 6-tuple $M \stackrel{\text{def}}{=} \langle \mathcal{S}, \mathcal{A}, T, r, \gamma, \mu_0 \rangle$ where $\mathcal{S}$ is the set of states, $\mathcal{A}$ the set of actions, $T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$ is the state transition function $T(s, a, s')$ defining the probability that dictates the evolution from $s \in \mathcal{S}$ to $s' \in \mathcal{S}$ after taking the action $a \in \mathcal{A}$, $r : \mathcal{S} \times \mathcal{A} \to [r_{min}, r_{max}]$, with $r_{max}, r_{min} \in \mathbb{R}$, is the reward function $r(s, a)$ that indicates what the agent gains when the system state is $s \in \mathcal{S}$ and action $a \in \mathcal{A}$ is applied, $\gamma \in [0, 1)$ is called the discount factor and $\mu_0 : \mathcal{S} \to [0, 1]$ is the distribution over initial states: $\sum_{s \in \mathcal{S}} \mu_0(s) = 1$.*

**Definition 2.** *A policy is a mapping from states to a probability distribution over actions, such as $\pi : \mathcal{A} \times \mathcal{S} \to [0, 1]$.*

**Definition 3.** *Solving an MDP amounts to finding a policy $\pi^*$ which, $\forall s \in \mathcal{S}$, maximizes the value function:*

$$V_M^\pi(s) \stackrel{\text{def}}{=} \mathbb{E}_{\substack{A_t \sim \pi \\ S_t \sim T}} \left[ \sum_{t=0}^{\infty} \gamma^t r(S_t, A_t) \,\middle|\, S_0 = s \right]. \tag{1}$$

It has been proved that an MDP for which the value function is defined as Eq. (1) admits a deterministic optimal policy (a map from states to actions) [36]:

$$\pi^*(s) = \operatorname{argmax}_\pi V_M^\pi(s). \tag{2}$$

**Definition 4.** *The performance of a policy $\pi$ in an MDP M with value function $V_M^\pi$ is defined as:*

$$u^\pi(M) = \mathbb{E}_{S \sim \mu_0}[V_M^\pi(S)]. \tag{3}$$

**Definition 5.** *A BMDP is a 8-tuple $\beta \stackrel{\text{def}}{=} \langle \mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \tau, \rho, \gamma, \mathcal{B} \rangle$ where $\mathcal{S}$ is the set of states; $\mathcal{A}$ the set of actions; $\mathcal{T}$ is a parametric family of transition functions $T$ of any MDP compatible with $\mathcal{S}$ and $\mathcal{A}$: $\mathcal{T} = \{T : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1] \text{ s.t. } \sum_{s' \in \mathcal{S}} T(s, a, s') = 1\}$; $\mathcal{R}$ is a parametric family of reward functions $r$ of any MDP compatible with $\mathcal{S}$ and $\mathcal{A}$: $\mathcal{R} = \{r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{min}, r_{max}]\}$; $\tau$ is a non-informative prior distribution uniform over $\mathcal{T}$: $\int_{T \in \mathcal{T}} d\tau = 1$ with $\tau \geq 0$; $\rho$ is a non-informative prior distribution uniform over $\mathcal{R}$: $\int_{r \in \mathcal{R}} d\rho = 1$ with $\rho \geq 0$; $\gamma \in [0, 1)$ is the discount factor, and $\mathcal{B} = \{(s_t, a_t, r_t, s_{t+1})\}$ is a batch of transitions generated by acting in a fixed, unknown MDP compatible with $\mathcal{S}$ and $\mathcal{A}$ and initial state distribution $\mu_0$.*

For instance, in a finite state and action spaces environment, $\mathcal{T}$ is the set of all $|\mathcal{S}| \times |\mathcal{A}|$ different discrete distributions and $\tau$ is made of $|\mathcal{S}| \times |\mathcal{A}|$ uniform (uninformative) Dirichlet probability density functions – the conjugate prior of the said distribution.

**Definition 6.** *$\tau_p$ is a posterior distribution over $\mathcal{T}$ obtained by updating the uniform (uninformative) Dirichlet prior $\tau$ with the information contained in $\mathcal{B}$.*

In particular, the $|\mathcal{S}|$ probability values $X_i$ with $i \in \{1, \ldots, \mathcal{S}\}$ describing the probability of $(S = s^*, A = a^*) \rightarrow (S' = s_i)$ can be distributed as:

$$\tau_p^{s^*, a^*} \left( x_1, \ldots, x_{|\mathcal{S}|} \big| n_1, \ldots, n_{|\mathcal{S}|} \right) = \Gamma(\nu) \prod_{i=1}^{|\mathcal{S}|} \frac{x_i^{n_i}}{\Gamma(n_i + 1)} \tag{4}$$

where, $\Gamma$ is the Euler gamma function, $n_i$ counts how many times the transition $(s^*, a^*) \rightarrow s_i$ appears in $\mathcal{B}$ and $\nu = \sum_{k=1}^{|\mathcal{S}|} (n_k + 1)$.

*Remark.* Notice that the mode (the most likely configuration) of the posterior in Eq. (4) is given by $\hat{x}_i = \frac{n_i}{\sum_{k=1}^{|\mathcal{S}|} n_k}$ while its expected value is $\mathbb{E}_{\tau_p}[X_i] = \frac{n_i + 1}{\nu}$.

Since we consider discrete environments, the most likely transition model with respect to $\tau_p$ is the one for which the transition probabilities are given by the transition frequencies in $\mathcal{B}$. We refer to these distributions as the trivial model, noted as $\hat{T}$. We note that a similar reasoning can be applied to $\mathcal{R}$ and $\rho$, however for simplicity's sake we assume to know the reward function $r$ in this work.

*Remark.* It would be possible to define a prior over the initial states and obtain a posterior taking into account the information contained in the batch $\mathcal{B}$. For simplicity, we will also assume that $\mu_0$ is known.

**Definition 7.** *A solution to a BMDP $\beta$ is a policy $\pi$ which maximizes the following utility function:*

$$\mathcal{U}_\beta^\pi \stackrel{\text{def}}{=} \mathbb{E}_{M \sim \tau_p} [u^\pi(M)] \tag{5}$$

*where, $u^\pi(M) \stackrel{\text{def}}{=} \mathbb{E}_{S \sim \mu_0} \left[ V_M^\pi(S) \right]$ is the expected value of an MDP, averaged on the initial state, with transition function sampled from $\tau_p$.*

The optimal performance with respect to Eq. (5) will be the one that, on average, works the best on the BMDP $\beta$ when the model is distributed according to the Bayesian posterior:

$$\mathcal{U}_\beta^* = \max_\pi \mathcal{U}_\beta^\pi \tag{6}$$
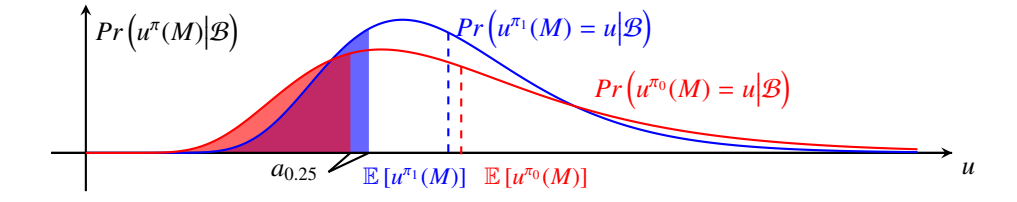
*Remark.* Since the true MDP is unknown, leveraging the Bayesian framework is an elegant way to incorporate uncertainty. However, the additional expected value makes Eq. (5) hard to be computed with Bellman's recursive approaches or approximated with temporal differences methods. Indeed, the objective stated in Eq. (6) is NP-hard [32].

## 3. Risk-aware measures

We advocate that solving a BMDP deals with uncertainty more elegantly than solving an MDP for the most likely model $\hat{M} = (\hat{T}, \hat{r})$. However, the utility function defined in Eq. (5) does not minimize the risk of obtaining a bad performant policy in the real environment.

For instance, let $Pr(u^\pi(M) = u|\mathcal{B})$ be the probability density function (pdf) of the performance of the policy $\pi$ when the model is distributed according to the Bayesian posterior $\tau_p$. Given two policies $\pi_0$ and $\pi_1$, we can have that $\mathbb{E}[u^{\pi_0}(M)] > \mathbb{E}[u^{\pi_1}(M)]$ (see Figure 1). In this case, following the BMDP optimization criterion, $\pi_0$ is better than $\pi_1$. However, when fixing a value $u$ less than both expected values, it can happen that $Pr\left(u^{\pi_0}(M) < u\middle|\mathcal{B}\right) > Pr\left(u^{\pi_1}(M) < u\middle|\mathcal{B}\right)$.

**Figure 1** Pdf for the performance $u^\pi(M) = u$ given the batch $\mathcal{B}$ and two different policies $\pi_0$ and $\pi_1$. Dashed lines correspond to the expected value. Both curves are filled up to their 0.25-quantiles ($a_{0.25}$). The red curve has a higher expected value while the blue one corresponds to a safer policy.



With this in mind, it can be useful to define risk-aware utility functions. Risk measures are widely studied in mathematical finance [37] since they are a way to rationally quantify risk. Their application to MDPs under model uncertainty was investigated in the work in [35].

In the following, we introduce two risk measures: the Value at Risk (*VaR*) and the Conditional Value at Risk (*CVaR*) both inspired by the work in [38].

Let $M$ be a random variable governed by a probability measure $Pr$ on its domain $\mathbb{M}$, and $u : \mathbb{M} \to \mathbb{R}$ be a measurable function such that $\mathbb{E}[u(M)] < +\infty$. Following the work in [38], we define the cumulative density function of $u(M)$ as:

$$\Psi(a) = Pr(u(M) \leq a) \tag{7}$$

While in finance or insurance industry, losses that should be minimized (by looking for the optimal decision) are considered, in the MDP framework, the function $u$ is a utility function that should be maximized (by seeking the optimal strategy). Let us consider a (low) risk level $q \in (0, 1)$, that corresponds to the (high) confidence level in [38].

**Definition 8** (Value at Risk). The Value at Risk (*VaR*) of the utility function $u$, at the risk level $q$ is

$$a_q = \min\{a \in \mathbb{R}|\Psi(a) \geq q\}. \tag{8}$$

The minimum in 8 is attained because $\Psi$ is non-decreasing and right continuous. The definition of the Conditional Value at Risk is slightly different from the one in [38], because again, in the MDP context, the lowest gains (to be maximized) are considered, and not the highest losses (to be minimized) like in finance.

**Definition 9** (Conditional Value at Risk). Let the cumulative density function of $u(M)$ conditional to $\{u(M) \leqslant a_q\}$ be

$$\Psi_q(a) = \begin{cases} \dfrac{\Psi(a)}{\Psi(a_q)} & \text{for } a \leq a_q, \\ 1 & \text{for } a > a_q. \end{cases} \tag{9}$$

The Conditional Value at Risk (*CVaR*) of the utility function $u$ at risk level $q$ is the expectation of the variable drawn by the cumulative density function 9

$$\phi_q = \mathbb{E}_{X \sim \Psi_q}[X]. \tag{10}$$

Note that $\Psi(a_q)$ may be higher than $q$, that is why $\Psi(a_q)$ is used instead of $q$ in the denominator of the Eq. 9, following the rescaling solution proposed in [38] in case of a probability atom at $a_q$. With these definitions, the (Conditional) Value at Risk of the utility function $u$ at risk level $q$ is also the (Conditional) Value at Risk of $u$ at risk level $\Psi(a_q) \geq q$.

*Remark.* Since the space of possible transition functions, that can be sampled from a Dirichlet posterior $\tau_p$, has the cardinality of the continuum, the distribution of the performance with respect to the uncertainty is continuous.

### 3.1. Risk-aware Solutions to BMDPs

In the following, the BMDP utility of Eq. (5) is generalized to take the risk into account.

**Definition 10.** *Let $\beta$ be a BMDP and let $V_M^\pi(s)$ be the value function at state $s$ while following a policy $\pi$ in the MDP $M$ with transitions distributed according to the posterior $\tau_p$. Let also $Pr(u^\pi(M)|\mathcal{B})$ be the pdf over the possible values assumed by $u^\pi(M) = \mathbb{E}_{S \sim \mu_0}\left[V_M^\pi(S)\right]$ given the batch $\mathcal{B}$. Then a risk-aware utility function is defined as:*

$$\mathcal{U}_{\beta,\sigma}^\pi \stackrel{\text{def}}{=} \sigma_{M \sim \tau_p}\left[u^\pi(M)\right] \tag{11}$$

*where, $\sigma$ is a risk measure.*

*Remark.* As a consequence of Definition 9 if $\sigma$ is the Conditional Value at Risk at risk level 1 then the BMDP utility of Eq. (5) is a particular case of Eq. (11).

## 4. Solving Offline a Risk-aware BMDP

The expectation over the distribution of models makes the solution of a BMDP an intractable computational task. Moreover, a Risk-aware BMDP also presents an additional difficulty: the risk measure requires an estimate of the quantiles of the unknown value distribution given a policy. An analytical maximization of the performance defined in Eq. (11) is often either impossible or too computationally demanding. In order to tackle the maximization problem, a valuable choice is resorting to a Monte Carlo estimate of the performance. We will then look for a sub-optimal policy, rather than an optimal one, by constraining the search to a set of candidate policies $\Pi$.

However, what number $L_\pi \in \mathbb{N}$ of models would be necessary to be sampled in order to have an accurate estimate of the performance of a policy within a chosen confidence interval? Ideally, $L_\pi$ should be as small as possible because Policy Evaluation has to be performed $L_\pi$ times in order to obtain the Bayesian posterior distribution of values assumed by the Value Functions.

Fortunately, this problem has been addressed by the work in [39]. It proposes a procedure that allows iteratively sampling values from a distribution whose quantile is required until the estimate of the said quantile will fall within a confidence interval with a required probabilistic significance.

In the present work, we exploit the idea of estimating a quantile through sampling to propose the Monte Carlo Confident Policy Selection (MC2PS) algorithm. MC2PS is presented in Algorithm 1. MC2PS identifies a robust policy for a Risk-aware Bayesian MDP among a set of candidate policies. In detail, for a given set of policies $\Pi$ and for every policy $\pi \in \Pi$, the algorithm

---

**Algorithm 1** MC2PS

**Input**: set of policies $\Pi$, significance level $\alpha \in [0, 1]$, sampling batch size $k \in \mathbb{N}$, relative error tolerance $\varepsilon_{rel} \in (0, 1]$, posterior distribution $\tau_p$, risk level $q \in (0, 1)$, risk measure $\sigma$, initial state distribution $\mu_0$, evaluation discount factor $\gamma$.

**Output**: best policy $\pi^*$.

1: **for** $\pi \in \Pi$ **do**
2:      $\mathcal{U}_{\beta,\sigma}^{\pi} \leftarrow \text{RiskEvaluation} \left( \pi, \sigma, \tau_p, \mu_0, \varepsilon_{rel}, \alpha, q, k, \gamma \right)$
3: **end for**
4: **return** $\pi^* = \arg\max_{\pi \in \Pi} \mathcal{U}_{\beta,\sigma}^{\pi}$
5:
6: **procedure** RiskEvaluation
7:      **Input**: policy $\pi$, risk measure $\sigma \in \{VaR, CVaR\}$, posterior distribution $\tau_p$, initial state distribution $\mu_0$, relative error threshold $\varepsilon_{rel} \in [0, 1]$, significance level $\alpha \in [0, 1]$, risk level $q \in (0, 1)$, sampling batch size $k \in \mathbb{N}$, evaluation discount factor $\gamma$.
8:      Initialize $u^{\pi} = \varnothing$
9:      *(the loop estimates the quantile needed in Eq. (11))*
10:      **repeat**
11:          **for** $j \in \{1, \dots, k\}$ **in parallel do**
12:              Sample $\mathcal{M}_j \sim \tau_p$
13:              $V_{\mathcal{M}_j}^{\pi}(s) \leftarrow$ Policy Evaluation on model $\mathcal{M}_j$
14:              $u^{\pi}(\mathcal{M}_j) \leftarrow \mathbb{E}_{S \sim \mu_0}[V_{\mathcal{M}_j}^{\pi}(S)]$ Eq. (3)
15:              $u^{\pi} \leftarrow$ **append** $u^{\pi}(\mathcal{M}_j)$
16:          **end for**
17:          $L_{\pi} \leftarrow |u^{\pi}|$
18:          Sort $u^{\pi}$ in increasing order
19:          Find $(g, h) \in \mathbb{N}^2$ such that $|h - g|$ is minimal and:
20:              $Pr(u_g^{\pi} \leq a_q < u_h^{\pi}) = \left( \sum_{i=g}^{h-1} \binom{L_{\pi}}{i} q^i (1-q)^{L_{\pi}-i} \right) > 1 - \alpha;$
21:      **until** :
22:          $u_h^{\pi} - u_g^{\pi} < \varepsilon_{rel} \cdot \left( u_{L_{\pi}}^{\pi} - u_1^{\pi} \right)$
23:      **if** $\sigma = VaR$ **then**
24:          $\widehat{a_q} \leftarrow u_g^{\pi}$
25:          **return** $\widehat{a_q}$
26:      **if** $\sigma = CVaR$ **then**
27:          $\widehat{\phi_q} \leftarrow \frac{1}{g} \sum_{i=1}^{g} u_i^{\pi}$
28:          **return** $\widehat{\phi_q}$
29: **end procedure**

---

incrementally samples $k$ transition models from $\tau_p$ and performs Policy Evaluation in parallel for each one of them until the stopping criterion is reached (see the RiskEvaluation procedure in Alg. 1). The stopping criterion guarantees that the estimate of the $q$-quantile is statistically well approximated with a significance level $\alpha$ within a dynamically sampled confidence interval whose width is smaller than $\varepsilon_{rel}$ (lines 19-22) given the total $L_\pi$ models sampled.

Indeed, being $u_i^\pi := u^\pi(M_i)$ the list of ordered performance values obtained from the sampled $L_\pi$ models with $i \in \{1, \ldots, L_\pi\}$, the probability that the elements with indices $h$ and $g$ of this list are bounding $a_q$, will be given by the probability of the union of all the (incompatible) events that lead to $u_g^\pi \leq a_q \leq u_h^\pi$ (see Figure 2). In detail, let $a_q$ be the theoretical Value at Risk of $u^\pi(M)$ at risk level $q$. Let us denote sampled utility values in increasing order by $u^\pi(M_1) \leq \ldots \leq u^\pi(M_{L_\pi})$, and suppose that the utility distribution has no probability atom at $a_q$: $\forall 1 \leq i \leq L_\pi, \Psi(a_q) := Pr(u^\pi(M_i) \leq a_q) = q$. Let us introduce the random variables

$$B_i = \mathbb{1}_{\{u^\pi(M_i) \leq a_q\}} = \begin{cases} 1, & \text{if } u^\pi(M_i) \leq a_q, \\ 0, & \text{otherwise.} \end{cases} \tag{12}$$

The random variables $B_i$ are drawn by a Bernoulli distribution with parameter $q$. The random variable $B = \sum_{i=1}^{L_\pi} B_i$ is the number of sampled utilities that are lower than $a_q$, drawn by a binomial distribution with parameters $L_\pi$ and $q$. The event $\{u^\pi(M_g) \leq a_q < u^\pi(M_h)\}$ is $\cup_{i=g}^{h-1}\{B = i\}$, *i.e.* the event "there are exactly $g, g + 1, \ldots,$ or $h - 1$ sampled utility values that are lower than $a_q$". Using the binomial distribution formula, the probability of this event is

$$Pr\Big(u^\pi(M_g) \leq a_q < u^\pi(M_h)\Big) = \sum_{i=g}^{h-1} Pr(B = i) = \sum_{i=g}^{h-1} \binom{L_\pi}{i} q^i (1 - q)^{L_\pi - i}. \tag{13}$$

Hence by imposing constraint $\sum_{i=g}^{h-1} \binom{L_\pi}{i} q^i (1 - q)^{L_\pi - i} > 1 - \alpha$ when selecting indices $r$ and $s$, we ensure that

$$Pr\Big(u^\pi(M_g) \leq a_q < u^\pi(M_h)\Big) > 1 - \alpha, \tag{14}$$

*i.e.* we get probabilistic bounds computed from the sampled utility values.

Note that, if there is a probability atom at $a_q$, *i.e.* $Pr(u^\pi(M_i) = a_q) > 0$, the previous reasoning cannot be applied directly in the case where $q < \Psi(a_q) := Pr(u^\pi(M_i) \leq a_q)$. However, we can write
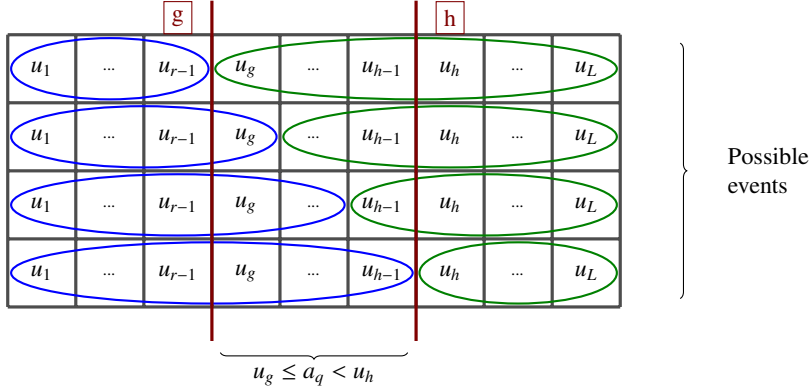
$$q_- := Pr(u^\pi(M_i) < a_q) \leq q < q_+ := Pr(u^\pi(M_i) \leq a_q), \tag{15}$$

and one can show that selected indices $g$ and $h$ are non decreasing with $q$. Thus, using the risk level $q$, selected indices are higher than those that would be selected using $q_-$, and lower than those that would be selected using $q_+$, both corresponding to utility values bounding the location of the probability atom, *i.e.* the Value at Risk $a_q$, with probability $1 - \alpha$.

Eventually, the algorithm leverages the estimate of both the Value at Risk and of the policy value achieved on sampled models to obtain an estimate of the utility function $\mathcal{U}_{\beta,\sigma}^\pi$ for a specific risk measure $\sigma$ and risk level $q$. For instance, return the estimate of $a_q$ if $\sigma = VaR$ or $\phi_q$ if $\sigma = CVaR$ (lines 23-28). Finally, once the utility function has been estimated for every policy, it outputs the one that maximizes it (line 4).

*Remark.* Let $\Lambda$ be the total number of models sampled to estimate the quantile of the performance distribution among policies: $\Lambda = \sum_{\pi \in \Pi} L_\pi$. MC2PS performs Policy Evaluation $\Lambda$ times. The size of the space of all applicable policies of a finite state and action space MDP is $|\Pi| = |\mathcal{A}|^{|\mathcal{S}|}$. It goes

**Figure 2** Example of how the estimate in Algorithm 1 works: imagine you have an ordered list with $L$ values $u_i$, $i \in \{1, \ldots, L\}$ represented in the figure as rows. The probability of the event in Eq. (14) is $\binom{L}{g} q^g (1-q)^{L-g}$ and corresponds to the probability of the random variable $B$ defined in Eq. (12) to assume all integer values between $g$ and $h-1$. The said probability is the sum of the probability of the events $B = i$ with $g \leq i < h$. In the figure, every addend is represented as a row. In blue are encircled the values of $u_i$ smaller than $a_q$ and in green the ones bigger. The algorithm looks for the indices $(g, h) = argmin_{(g,h)}|g - h|$ such that $Pr(u_g \leq a_q < u_h) > 1 - \alpha$ and $u_h - u_g < \varepsilon$ where $\varepsilon$ is an error term dictating the maximum acceptable size of the confidence interval.



without saying that looking over the whole policy space can be practically intractable even for not-so-big MDPs. Nevertheless, restricting the research to a subset of policies could be a viable solution also for big MDPs, also considering that the Policy Evaluations are carried out in parallel.

### 4.1. Exploitation vs Caution (EvC)

Reference [25] shows that the policy obtained by solving an MDP $\hat{M} = (\mathcal{S}, \mathcal{A}, \hat{T}, r, \gamma^*)$, trivially learned from a batch of experiences collected from another MDP $M = (\mathcal{S}, \mathcal{A}, T, r, \gamma_{ev})$, with $\gamma^*$ a discount factor such as $\gamma^* \leq \gamma_{ev}$, is more efficient in $M$ than the policy obtained by solving $\hat{M}$ using $\gamma_{ev}$. The reason is that $\hat{T}$ is an approximation of $T$ and may not be trusted for long-term planning horizons. The selection of the best $\gamma^*$ optimizes a trade-off between the exploitation of the information contained in the batch and the necessity of being cautious since the model estimate is not perfect.

Inspired by the conclusions of reference [25], and also guided by the intuition that the model $M$ that generated the data will be different from $\hat{M}$, but hopefully *close* to it, we expect that the policies obtained by solving another MDP $\tilde{M} = (\mathcal{S}, \mathcal{A}, \tilde{T}, r, \gamma)$ with $\tilde{T}$ *close* to $\hat{T}$ and $\gamma \leq \gamma_{ev}$ can be viable solutions for the Risk-aware Bayesian MDP.

Henceforth, the Exploitation *versus* Caution (EvC) algorithm is presented and schematized in Algorithm 2. EvC will search for a promising risk-aware policy by focusing the search on a set of candidate policies $\Pi$ computed with several baseline algorithms $\mathbb{A}$, which is further enriched by solving different MDPs $\tilde{M}$ and $\gamma \leq \gamma_{ev}$ values. Remember that the goal is to find a policy that is performant in the model $M$ taking into account that the agent does not have access to the probability values that define the model, but only to the pre-collected batch. Model uncertainty is framed within a Bayesian MDP. As already stated, we do not aim to find the optimal solution

---

**Algorithm 2** EvC

---

    **Input**: risk level $q \in [0, 1]$, significance level $\alpha \in [0, 1]$, sampling batch size $k \in \mathbb{N}$, relative error tolerance $\varepsilon_{rel} \in [0, 1]$, posterior distribution $\tau_p$, risk measure $\sigma$, initial state distribution $\mu_0$, set of discount factors $G$, number of models to solve $l \in \mathbb{N}$, $\mathcal{B}$ batch of transitions, evaluation discount factor $\gamma_{ev}$

    **Output**: best policy $\pi^*$.

1:  $\Pi \leftarrow$ GENERATEPOLICIES $\left( \tau_p, G, l \right)$

2:  **return** $\pi^* = $ MC2PS $\left( q, \alpha, k, \varepsilon_{rel}, \tau_p, \Pi, \sigma, \mu_0, \gamma_{ev} \right)$

3:

4:  **procedure** GENERATEPOLICIES

5:     **Input**: posterior distribution $\tau_p$, set of discount factors $G$, number of models $l \in \mathbb{N}$ to be solved, $\mathcal{B}$ batch of transitions.

6:     Initialize $\mathbb{M} = \left\{ l \text{ transition models} \sim \tau_p \right\} \cup \{\hat{T}\}$

7:     Initialize $\Pi = \varnothing$ (an empty set)

8:     Initialize $\mathbb{A} = \{SPIBB, BOPAH, BCR, NORBU\}$ (examples of baseline algorithms)

9:     **for** $(\gamma \in G, T \in \mathbb{M})$ **do**

10:       $\pi_{(T,\gamma)} = $ solution to the MDP with $T$ and $\gamma$

11:       Append $\pi_{(T,\gamma)}$ to $\Pi$ if $\pi_{(T,\gamma)} \notin \Pi$

12:     **end for**

13:     **for** *algorithm* $\in \mathbb{A}$ **do**

14:       $\pi_{algorithm} = $ solution to the Offline MDP with $\mathcal{B}$ and *algorithm*

15:       Append $\pi_{algorithm}$ to $\Pi$ if $\pi_{algorithm} \notin \Pi$

16:     **end for**

17:     **return** $\Pi$

18: **end procedure**

---

to the BMDP, but rather to select the best policy in terms of robustness within the ones in the candidate set.

In detail, EvC first generates candidate policies that will constitute the set $\Pi$ (line 1 calls GENERATEPOLICIES procedure). For this, starting from the batch the problem is solved using a portfolio constituted of state-of-the-art algorithms (line 8). On top of that the trivial MDP[1] $\hat{M}$ and $l$ additional MDPs are sampled from the Bayesian posterior $\tau_p$ obtained from the batch (line 6), and then solved with different values of $\{\gamma \in G | \gamma \leq \gamma_{ev}\}$ (lines 9-12). Recalling that $\gamma_{ev}$ is the discount factor of the Risk-aware Bayesian MDP. Note that the obtained set $\Pi$ has unrepeated solutions (line 11 and line 15). As a last step, MC2PS is launched with the obtained set of candidate policies $\Pi$ returning the best risk-aware solution $\pi^* \in \Pi$ (line 2).

*Remark.* Note, if we test over 9 different discount factors, such as $G = \{0.1, 0.2, \ldots, \gamma_{ev} = 0.9\}$, and 5 different $(l = 5)$ MDPs $\tilde{M}$ (including $\hat{M}$), then we solve $|\Pi| \leq 9l = 45$ MDPs to enrich the set of candidate policies within this approach.

---

[1]The trivial MDP $\hat{M}$ is a straightforward MDP estimate using the batch $\mathcal{B}$. For instance, in the case of a discrete MDP this is equivalent to the model that maximizes the likelihood of $\mathcal{B}$, *i.e.* the one whose transition probabilities are obtained from the frequencies of transitions in the batch.

## 4.2. Theoretical guarantees

Since EvC searches for the policy $\pi \in \Pi$ that maximizes the criterion of Eq. (11), the Algorithm 2, rather than yielding a sub-optimal solution to the Risk-aware BMDP, can be seen as a policy selection approach. Assuming that the Bayesian posterior $\tau_p$ efficiently encodes the model uncertainty, EvC outputs a policy whose performance in the real environment is guaranteed in probability to be greater than some value that changes with respect to the chosen risk-aware measure. In a simpler way, we can provide theoretical guarantees on the estimate of the quantile needed to compute the risk-aware utility function that will be eventually maximized over the set of candidate policies.

**Theorem 1.** *Let $\pi \in \Pi$ be a candidate policy and $u^\pi(M_g)$ be an estimate of the Value at Risk of $u^\pi(M)$ at risk level $q$ calculated through EvC. Let $u^\pi(M)$ be the performance of $\pi$ with $M$ distributed according to the Bayesian posterior $\tau_p$. The performance of $\pi$ in this MDP $M$ is greater than the estimate of $a_q$ with probability:*

$$Pr(u^\pi(M) \geq u^\pi(M_g)) \geq (1-q)(1-\alpha). \tag{16}$$

*Proof.* Note that $\{u^\pi(M) \geq a_q\} \cap \{a_q \geq u^\pi(M_g)\} \subseteq \{u^\pi(M) \geq u^\pi(M_g)\}$, where $a_q$ denotes the Value at Risk of $u^\pi(M)$ at risk level $q$. The two events of the intersection respectively depend on two independent random variables – a future performance $u^\pi(M)$, that could be obtained by acting according to the policy $\pi$, and a Value at Risk estimate $u^\pi(M_g)$, whose randomness is the result of the sampling procedure in the Algorithm 1. The previous inclusion allows writing $Pr(u^\pi(M) \geq u^\pi(M_g)) \geq Pr(u^\pi(M) \geq a_q) \cdot Pr(a_q \geq u^\pi(M_r)) \geq (1-q) \cdot Pr(a_q \geq u^\pi(M_g)) \geq (1-q)(1-\alpha)$. The last inequality is ensured by the quantile estimation (lines 19-22 in Algorithm 1), and the previous one by the definition of $a_q$. Therefore, we get the equation 16. $\qquad\square$

*Remark.* When the risk-aware measure used in EvC is *VaR* the lower bound on $u^\pi(M)$ ($u^\pi(M_g)$) in the Proof of Theorem 1 is maximized over the policies. If the risk-aware measure is *CVaR* the empirical expected value over the $q$-fraction of low-performant policies is maximized.

*Remark.* Since $u^\pi(M_1) \leq u^\pi(M_2) \leq \ldots \leq u^\pi(M_g)$, then $\frac{1}{g}\sum_{i=1}^g u^\pi(M_i) \leq u^\pi(M_g)$, therefore the same lower bound in probability is also valid for the *CVaR* utility function:

$$Pr\left(u^\pi(M) \geq \frac{1}{g}\sum_{i=1}^g u^\pi(M_i)\right) \geq Pr(u^\pi(M) \geq u^\pi(M_g)) \geq (1-q)(1-\alpha). \tag{17}$$

Note that the sampling procedure in the Algorithm 1 ensures that $Pr(|u^\pi(M_g) - a_q| \geq \varepsilon) \leq \alpha$, with $\varepsilon = \varepsilon_{rel} \cdot (u^\pi(M_{L_\pi}) - u^\pi(M_1))$. If a practitioner wants such a probabilistic bound on the precision of the estimate of $\phi_q$, she/he should sample additional models $N \in \mathbb{M}$ from the posterior, to select $n$ independent models such that $u^\pi(N) \leq u^\pi(M_g)$, where $u^\pi(M_g)$ is given by the sampling procedure of Algorithm 1. The new estimate of $\phi_q$ computed from these $n$ new models benefits from the following theorem.

**Theorem 2.** *Let $\pi \in \Pi$ be a candidate policy, $N_i \in \mathbb{M}$ be one of the n new sampled models from the posterior $\tau_p$ such that $\forall i, u^\pi(N_i) \leq u^\pi(M_g)$, with $u^\pi(M_g)$ calculated through EvC, and $\overline{U} = \frac{1}{n}\sum_{i=1}^n u^\pi(N_i)$ be the new estimate of $\phi_q$. This new estimate of the Conditional Value at Risk of $u^\pi$ at risk level q respects the following inequality:*

$$Pr\left(|\overline{U} - \phi_q| \geq t\right) \leq 2\exp\left(-\frac{2nt^2}{(u^\pi(M_g) - \xi)^2}\right) + \alpha, \tag{18}$$

with $\xi = \inf_{m \in \mathbb{M}} u^\pi(m)$, or any other lower bound of $u^\pi$ as, for instance, $\frac{r_{min}}{1-\gamma}$. Note that $\xi \geq 0$ if the reward values are known to be non-negative.

*Proof.* By using the law of total probability, and upper bounding some probability values by 1,

$$
\begin{aligned}
Pr\big(|\overline{U} - \phi_q| \geq t\big) &= Pr\big(|\overline{U} - \phi_q| \geq t \,\big|\, \forall i, u^\pi(N_i) \leq a_q\big) Pr(\forall i, u^\pi(N_i) \leq a_q) \\
&\quad + Pr\big(|\overline{U} - \phi_q| \geq t \,\big|\, \exists i \text{ s.t. } u^\pi(N_i) > a_q\big) Pr(\exists i \text{ s.t. } u^\pi(N_i) > a_q) \\
&\leq Pr\big(|\overline{U} - \phi_q| \geq t \,\big|\, \forall i, u^\pi(N_i) \leq a_q\big) + Pr(\exists i \text{ s.t. } u^\pi(N_i) > a_q).
\end{aligned}
$$

The probability value on the right is lower than $Pr(u^\pi(M_g) > a_q) \leq Pr\big(a_q \notin [u^\pi(M_g), u^\pi(M_h))\big) \leq \alpha$ using the inequality $Pr(u^\pi(M_g) \leq a_q) > 1 - \alpha$ from lines 19-22 of Algorithm 1. What follows only depends on the definition of $\phi_q$ as the expected value up to $a_q$, and Hoeffding's inequality. $\quad\square$

### 4.3. Consequences and applications

The purpose of Offline Learning is that of providing behavioral policies to be applied by real-world automated agents. Thus reducing the risk at the expense of a longer computational phase is not only commendable but compulsory. Will the policy obtained through MC2PS and EvC be *good* or entirely-risk free? This goes beyond the theoretical guarantees provided by the algorithms since its outputs depend not only on the characteristics of the environment and on the set of candidate policies but also on the quality and variety of the batch. A batch of transitions that is too small or too concentrated in the same region of the state-action space may result in policies that, even if they are guaranteed to handle the risk better than the trivial one, can still be catastrophic.
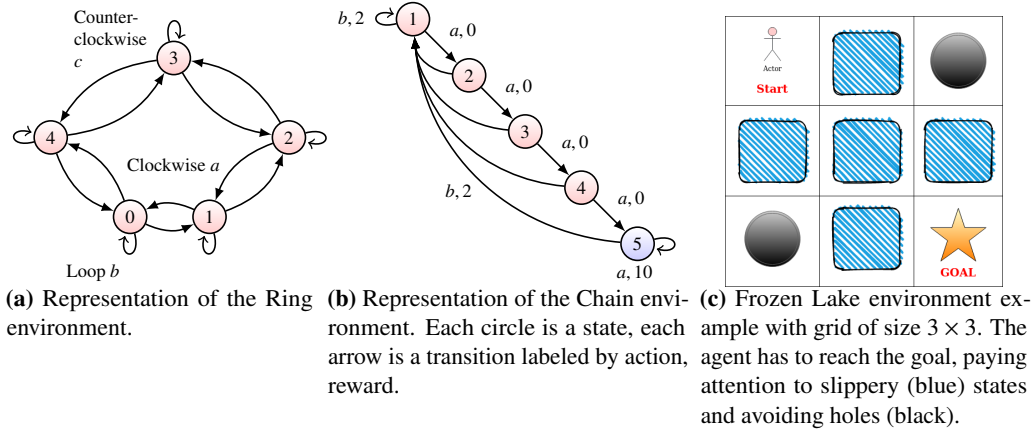
## 5. Experiments

In order to evaluate the proposed approach, we selected three small MDPs and hence easy-to-study stochastic environments endowed with diversified characteristics: two planning environments without absorbing states, Ring (5 states, 3 actions), and Chain (5 states, 2 actions). The former consisting in the stabilization of the agent in a particular non-absorbing goal with stochastic drift and the latter presenting cycles; and the Random Frozen Lake (RFL) environment, a re-adaptation of Frozen Lake from Open AI Gym suite [40] ($8 \times 8$ grid world with fatal absorbing states).

### 5.1. Environments' description

*Ring.* This environment is described by five states: $\{0, \ldots, 4\}$, forming a single loop. Three actions are possible: $a$, $b$, and $c$. The agent starts in state 0. The action $a$ will move it to the state $s - 1$ with probability 1.0 (e.g. when in 4 it moves to 3) if $s = 0, 1, 3$, and with probability 0.5 if it is elsewhere. With the action $b$ the agent will remain in the same state with probability 0.8 and move to the left or to the right with probability 0.1 if it is in state $s = 0, 1, 3$, if it is in state 2 or 4 it will move with probability 1. The action $c$ will move the agent to the right with probability 0.9 and it will not move with probability 0.1 if it is in state $s = 0, 1, 3$. Otherwise, the same effects will apply, but with probability 0.5. The agent earns an immediate reward $r = 0.5$ if it moves from $2 \rightarrow 3$ or $4 \rightarrow 3$, and $r = 1$ for any transition $3 \rightarrow 3$. Elsewhere $r = 0$. A graphical representation is shown in Figure 3a.

**Figure 3** Environments illustration.



(a) Representation of the Ring environment.

(b) Representation of the Chain environment. Each circle is a state, each arrow is a transition labeled by action, reward.

(c) Frozen Lake environment example with grid of size $3 \times 3$. The agent has to reach the goal, paying attention to slippery (blue) states and avoiding holes (black).

*Chain.* This environment was proposed in the research in [28] and was adapted to the present study. There are five states with the topology of an open chain and two actions $a$ and $b$. The agent starts from the state most to the left. With action $a$ the agent moves to the right and receives an immediate reward $r = 0$ with probability 0.8. Once the agent is in the rightmost state, performing the first action lets him stay there and receive a reward $r = 10$ with probability 0.8. It slips back to the origin earning a reward $r = 2$ with probability 0.2. Action $b$ moves the agent to the origin state with probability 0.8 receiving a reward $r = 2$ or letting it go right with probability 0.2 earning $r = 0$. The optimal policy consists of applying action $b$ in the first state and action $a$ in the others. A representation is shown in Figure 3b.

*Random Frozen Lake (RFL).* The Frozen Lake Environment of the Open AI Gym suite [40] was edited for this study. The agent moves in a grid world ($8 \times 8$). It starts in the utmost left corner and it must reach a distant absorbing goal state that yields a reward $r = 1$. In the grid there are some holes. If it falls into a hole it is blocked there and it can not move anymore, obtaining from that moment an immediate reward $r = 0$. Unfortunately, the field is covered with ice and hence it is slippery. When the agent wants to move towards a nearby state it can slip with fixed probability $p$ and ends up in an unintended place. The grid is generated randomly assuring that there always exists a hole-free path connecting the start and the goal. Moreover, to each couple of action and non-terminal state $(a, s)$ is assigned a different immediate reward $r$ sampled at random between $(0, 0.8)$ at the moment of the generation of the MDP problem. The MDP itself does not have a stochastic reward, but the map and the rewards are randomly generated. A graphical representation (for a $3 \times 3$ grid) is shown in Figure 3c.

## 5.2. Setup

Given $(n, m) \in \mathbf{N}^2$, $m$ trajectories, with $n$ steps each, are generated following a random policy in each environment. We opted for a random data collecting procedure because we imagine using EvC in a scenario where both the developers and the autonomous agent are completely agnostic about the model dynamics and have no prior knowledge.

The true environment is assumed to be known for the a posteriori evaluation. The most likely transition model is inferred from the batch. The trivial MDP was then solved with the Policy

Table 1: Parameters and hyperparameters used during the simulations: $n$ is the number of steps in each trajectory contained in a batch; $l$ is the number of different models sampled from the prior in EvC (Algorithm 2); $\{N_\wedge\}$ is the set of different thresholds used in SPIBB; **fold** and **DOF** are the fold and degree of freedom hyper-parameters used in BOPAH; $\lambda$ is the soft robust hyper-parameter of NORBU. Bold values are displayed in the plots.

| Environment | $n$ | $l$ | $\{N_\wedge\}$ | fold | DOF | $\lambda$ |
|---|---|---|---|---|---|---|
| Ring | 8 | 3 | $\{\mathbf{1}, 2, 3, 5, 7, 10, 20\}$ | 2 | 20 | 0.5 |
| Chain | 8 | 3 | $\{\mathbf{1}, 2, 3, 5, 7, 10, 20\}$ | 2 | 20 | 0.5 |
| RFL ($8 \times 8$) | 15 | 10 | $\{\mathbf{1}, 2, 3, 5, 7, 10, 20\}$ | 2 | 20 | 0.5 |

Iteration algorithm and its relative performance in the true environment is obtained by Policy Evaluation. EvC data was computed with $\phi_{0.25}$ and $a_{0.25}$ (the first quartile). For each of these risk-aware measures, the following parameters (see Algorithm 2) were used: the set of discount factors $G = \{0.2, 0.4, 0.6, 0.8, 0.9\}$, the significance level $\alpha = 0.01$, the relative tolerance error $\varepsilon_{rel} = 0.01$, and the number $l$ of different models sampled from the prior is given in Table 1.

In the experiments, for a given batch size $N = nm \in \mathbf{N}$, 50 different batches were generated containing fixed size trajectories. The trajectory sizes used are also given in Table 1.

The chosen state-of-the-art algorithms that provide the base for the set of candidate policies are the following:

1. *Deterministic policies:* output by the following baselines[1], please notice that the quantile used for the robust and soft robust objectives in the algorithms is the same provided as general input for the estimate of EvC: **BCR** [22], **NORBU - Soft Robust *CVaR*** [24] (soft robust hyperparameter $\lambda = 0.5$);

2. *Stochastic policies:* output by the following algorithms [2]: **SPIBB** [3] receiving as input the batch collector policy, and, **BOPAH** [7] receiving as input the batch collector policy.

In our implementation of these baselines we only used intuitively tunable parameters (e.g. the discount factor).

*Remark.* We did not use MOPO [9] and MOReL [10] since: (1) they have usually been tested on continuous state MDPs driven by deterministic dynamics, while here we are tackling non-deterministic environments; (2) they highly rely on hyperparameter domain-dependent fine tuning which we did not do to fulfill the offline learning obligation.

In the evaluation phase, the discount factor is defined as $\gamma_{ev} = 0.9$. The others simulation parameters are provided in Table 1. Eventually, we also compared EvC with UnO by performing the risk-sensitive off-policy evaluation with UnO over the same set of candidate policies provided to EvC and then selected the one that maximized the risk-sensitive objective. While it is true that UnO, as other Importance Sampling based off-policy evaluation methods, should not be able to accurately evaluate deterministic policy, we still compare our approach to it because there are no other risk-sensitive off-policy evaluation approaches of our knowledge.

*5.3. Metrics*

We report metrics about the performance differences $\Delta U = u_{\beta,\sigma}^{\pi} - u_{\beta,\sigma}^{\pi_{trivial}}$ of the policies $\pi$ obtained with a specific algorithm (Eq. (3) using the utility function defined in Eq. (11)) and the

---

[1]The code was taken from the authors' Github repository: https://github.com/marekpetrik/craam2/tree/master/examples/evaluation/algorithms and readapted.

[2]The code was taken from the Github repository: https://github.com/KAIST-AILab/BOPAH and readapted.

performance obtained by solving the trivial model in the same setting and using the same batch of trajectories. This last value is normalized by the performance of the optimal policy. In particular, we consider: (1) the maximal $\Delta U$ obtained, (2) the mean value over all the different simulations, (3) the median over all simulations, and (4) the minimal $\Delta U$. The selected metrics provide insight into the validity of the approaches. We consider only the extrema of the distributions of the results (min, max), their median, and mean values since trying to estimate the whole distributions, and hence their quantiles could result in wrong conclusions if we are not sampling enough batches. For instance, in order to correctly estimate the Value at Risk at risk level $q = 0.25$ with a $\alpha = 0.01$ significance usually tens of thousands of samples are required. However, we are performing only hundreds of simulations with a fixed batch size $N$, which is enough for the selected metrics but definitely insufficient for the study of the whole distribution.

*Remark.* Please notice that the distribution whose statistics are displayed in the tables is not the one used to maximize Eq. (11) since it is a distribution over different starting batches collected with the same random policy and not the distribution that encodes the model uncertainty using the same starting dataset. Indeed, from a bayesian point of view the results are distributed along:

$$Pr\left(u^{\pi}(M), \mathcal{B}|\pi_{random}\right) = Pr\left(u^{\pi}(M)|\tau_p\right) Pr\left(\tau_p|\mathcal{B}\right) Pr\left(\mathcal{B}|\pi_{random}\right), \tag{19}$$

that represents the probability of collecting a batch $\mathcal{B}$ by collecting transitions using a random policy $\pi_{random}$ and hence observing the performance $u^{\pi}(M)$ by deploying a policy $\pi$. Note that there is a deterministic mapping among the posterior $\tau_p$ and the batch, therefore $Pr\left(\tau_p|\mathcal{B}\right)$ is a delta function.

### 5.4. Results and Discussion

For Ring and Chain, the results averaged over 100 different batches for each batch size $N \in \{8, 16, 24, 32, 40, 48, 56\}$ are displayed in Table 2. While for RFL the results averaged over 50 different batches for every batch size $N \in \{15, 30, 45, 60, 75, 90, 105\}$ are reported in Table 3.

Even if the datasets are composed of relatively short trajectories (Ring and Chain $n = 8$ time steps each, Random Frozen Lake $n = 15$ time steps each) in most cases UnO does not manage to evaluate the deterministic policies. Please note that UnO computes the Importance Sampling ratio for a trajectory $h$, a policy $\pi$ and a behavioral policy $\pi_\beta$ as

$$\rho_h^{\pi} = \prod_{i=1}^{n_h} \frac{\pi\left(s_i, a_i\right)}{\pi_\beta\left(s_i, a_i\right)} \tag{20}$$

where $n_h$ is the number of time steps of the trajectory $h$. However, this formulation assumes that both $\pi$ and $\pi_\beta$ are stochastic. In our formulation $\pi_\beta(s, a) = |\mathcal{A}|^{-1} \ \forall (s, a) \in \mathcal{S} \times \mathcal{A}$, but $\pi$ is stochastic only when it is the output of SPIBB or BOPAH. When $\pi$ is deterministic, the former equation can be rewritten as

$$\rho_h^{\pi} = \prod_{i=1}^{n_h} \frac{\delta_{\pi(s_i), a_i}}{\pi_\beta\left(s_i, a_i\right)} = |\mathcal{A}|^{n_h} \prod_{i=1}^{n_h} \delta_{\pi(s_i), a_i}. \tag{21}$$

This means that $\rho_h^{\pi} = |\mathcal{A}|^{n_h}$ if and only if all sequence of actions and states is consistent with the deterministic policy $\pi$, otherwise $\rho_h^{\pi} = 0$. It goes without saying that the probability that the ratio will be zero grows as a power of $|A|$ and exponentially in $n_h$. In particular, the probability that a sequence will be generated by the deterministic policy is in Ring $|A|^{-n_h} = 3^{-8} \approx 1.5 \times 10^{-4}$, in

Chain $2^{-8} \approx 3.9 \times 10^{-3}$ and in RFL $4^{-15} \approx 9.3 \times 10^{-10}$. Therefore, almost always UnO will pick a policy among SPIBB and BOPAH since the Importance Sampling ratio will be zero for other policies. If even the output of SPIBB and BOPAH will result in a zero Importance Sampling ratio, then the first policy in the candidate set (the trivial policy) will be picked. The said phenomenon is what happens most of the time. Therefore UnO alternates among the Trivial Policy, one among SPIBB and BOPAH, and once in a while it selects another approach.

*Ring.* Using $q = 0.25$ the best method according to the *Max*, *Mean* and the *Median* is NORBU with the *CVaR* Soft Robust objective (see Table 2). However, the most robust baseline in terms of worst-case performance is BOPAH. The distributions of results are asymmetric around $\Delta U$. In the cases of BCR and NORBU the *Mean* and the *Median* are approximately zero. In the cases of SPIBB and BOPAH the *Median* and the *Mean* are less than zero. Regarding the off-policy evaluation and selection methods, EvC with the *VaR* is the most performing one with respect to all the considered metrics.

*Chain.* In this environment every baseline except for SPIBB works the same with SPIBB being the worst in terms of *Min* (see Table 2). Regarding the off-policy evaluation and selection, all algorithms perform well since there is not really a substantial difference between the approaches (except for SPIBB).

*Random Frozen Lake (RFL).* We test the approaches in 4 different RFLs. The best approach in terms of overall metrics in 3 environments over 4 is again NORBU with the Soft Robust *CVaR* (see Table 3). SPIBB is the best in Environment 4. The best selection method is EvC with *VaR*/*CVaR* ($a_q$ / $\phi_q$) which provides identical performances.

Table 2: Statistics of the normalized performance difference $\Delta U$ between the reported algorithm (risk level $q = 0.25$) and the trivial policy averaged over batch size $N \in \{8, 16, 24, 32, 40, 48, 56\}$ with 100 different batches for size in Ring and Chain. On the right $\Delta U$ with the algorithm selected by EvC and UnO with $a_{0.25}$ and $\phi_{0.25}$. Notice that both EvC and UnO can pick also a policy obtained with a model solved with a different discount factor.

| Environment | Metrics | Baseline | | | | Selection Method | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | SPIBB | BOPAH | BCR | NORBU | $EvC_{a_{0.25}}$ | $EvC_{\phi_{0.25}}$ | $UnO_{a_{0.25}}$ | $UnO_{\phi_{0.25}}$ |
| | *Max* | 0.61 | 0.48 | 0.74 | **0.84** | **0.82** | 0.71 | **0.82** | 0.72 |
| Ring | *Mean* | -0.29 | -0.28 | -0.01 | **0.03** | **0.01** | -0.04 | -0.26 | -0.27 |
| | *Median* | -0.31 | -0.34 | **0.0** | **0.0** | **0.0** | **0.0** | -0.27 | -0.33 |
| | *Min* | -0.78 | **-0.68** | -0.82 | -0.71 | -0.82 | -0.82 | -0.96 | -0.96 |
| | *Max* | **0.55** | 0.54 | **0.55** | **0.55** | **0.55** | **0.55** | 0.54 | 0.54 |
| Chain | *Mean* | 0.0 | 0.01 | 0.01 | **0.02** | **0.01** | **0.01** | **0.01** | **0.01** |
| | *Median* | **-0.01** | **-0.01** | **-0.01** | **-0.01** | **-0.01** | **-0.01** | **-0.01** | **-0.01** |
| | *Min* | -0.38 | -0.16 | **-0.15** | **-0.15** | **-0.16** | **-0.16** | **-0.16** | **-0.16** |

In the following, we comment on the results obtained with EvC. Note this algorithm selects the policy that optimizes the (Conditional) Value at Risk over the first quartile ($q = 0.25$) starting from the set of candidate policies discussed in the last section.

In terms of risk awareness, after a global study over different batch sizes, EvC does not select the policy that produces the best values with respect to the considered metrics. Nevertheless, the policy selected by EvC is between the more robust ones. These results are shown in Figures 4, 5, and 6. In particular, our approach tends to opt for a policy from the ones obtained by solving several models with different discount factors $\gamma$ when the batch is small. The number of times such a policy is selected decreases to the benefit of 1) the trivial policy when the batch size $N$

Table 3: Statistics of the normalized performance difference $\Delta U$ between the reported algorithm (quantile order used $q = 0.25$) and the trivial policy averaged over batch size $N \in \{15, 30, 45, 60, 75, 90, 105, 120, 135\}$ with 50 different batches for size in different Random Frozen Lake environments.

| Environment | Metrics | Baseline | | | | Selection Method | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | SPIBB | BOPAH | BCR | NORBU | $EvC_{a_{0.25}}$ | $EvC_{\phi_{0.25}}$ | $UnO_{a_{0.25}}$ | $UnO_{\phi_{0.25}}$ |
| RFL Env. 1 | *Max* | **0.32** | 0.31 | 0.31 | **0.32** | 0.3 | 0.32 | **0.37** | 0.32 |
| | *Mean* | **0.05** | -0.04 | -0.04 | **0.05** | **0.05** | **0.05** | -0.02 | -0.05 |
| | *Median* | **0.04** | -0.07 | -0.04 | **0.04** | **0.04** | **0.04** | -0.01 | -0.08 |
| | *Min* | -0.25 | **-0.22** | -0.39 | -0.33 | -0.33 | -0.33 | -0.31 | **-0.22** |
| RFL Env. 2 | *Max* | 0.33 | 0.22 | 0.3 | **0.34** | **0.34** | **0.34** | 0.28 | 0.18 |
| | *Mean* | 0.02 | -0.07 | -0.05 | **0.06** | **0.06** | **0.06** | 0.0 | -0.07 |
| | *Median* | 0.01 | -0.08 | -0.06 | **0.06** | **0.06** | **0.06** | -0.01 | -0.08 |
| | *Min* | -0.21 | -0.22 | -0.29 | **-0.12** | **-0.12** | **-0.12** | -0.28 | -0.26 |
| RFL Env. 3 | *Max* | 0.3 | 0.23 | **0.43** | 0.36 | **0.36** | **0.36** | 0.35 | 0.18 |
| | *Mean* | 0.01 | -0.08 | 0.0 | 0.04 | 0.04 | 0.04 | -0.03 | -0.08 |
| | *Median* | -0.0 | -0.09 | 0.01 | **0.02** | **0.02** | **0.02** | -0.03 | -0.09 |
| | *Min* | **-0.16** | -0.3 | -0.36 | -0.27 | -0.27 | -0.27 | -0.29 | **-0.26** |
| RFL Env. 4 | *Max* | 0.32 | 0.22 | **0.36** | 0.31 | **0.31** | **0.31** | 0.27 | 0.22 |
| | *Mean* | 0.02 | -0.06 | 0.01 | **0.05** | **0.05** | **0.05** | -0.05 | -0.06 |
| | *Median* | 0.02 | -0.06 | 0.01 | **0.05** | **0.05** | **0.05** | -0.05 | -0.06 |
| | *Min* | -0.32 | -0.3 | -0.4 | **-0.29** | **-0.29** | **-0.29** | -0.39 | -0.3 |

**Figure 4** Policy selection rate by EvC $a_{0.25}$ and EvC $\phi_{0.25}$ in Ring for different batch sizes.
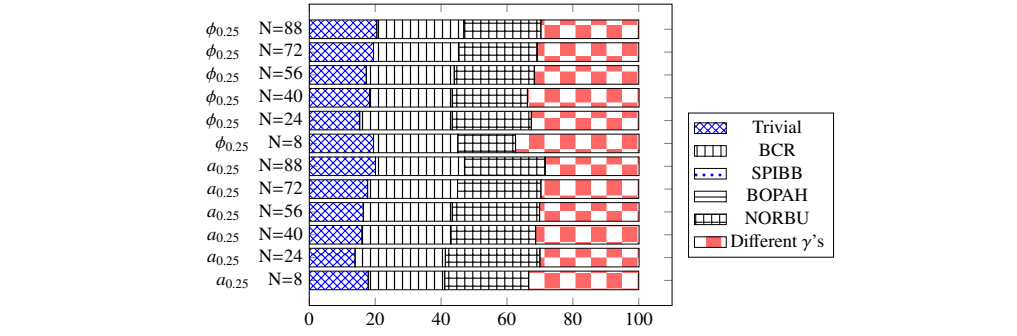


**Figure 5** Policy selection rate by EvC $a_{0.25}$ and EvC $\phi_{0.25}$ in Chain for different batch sizes.
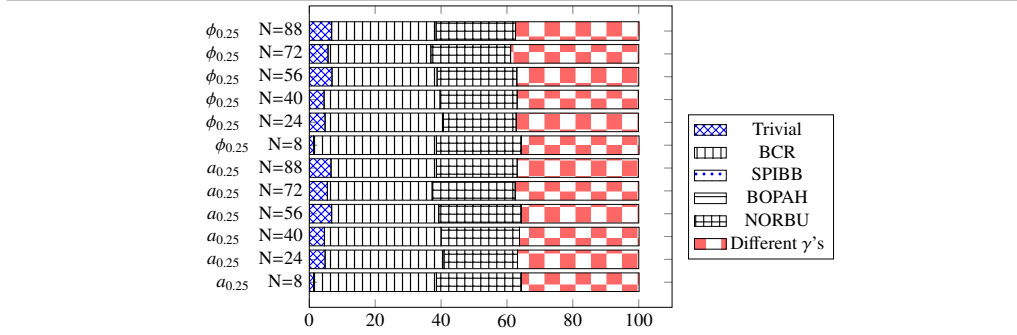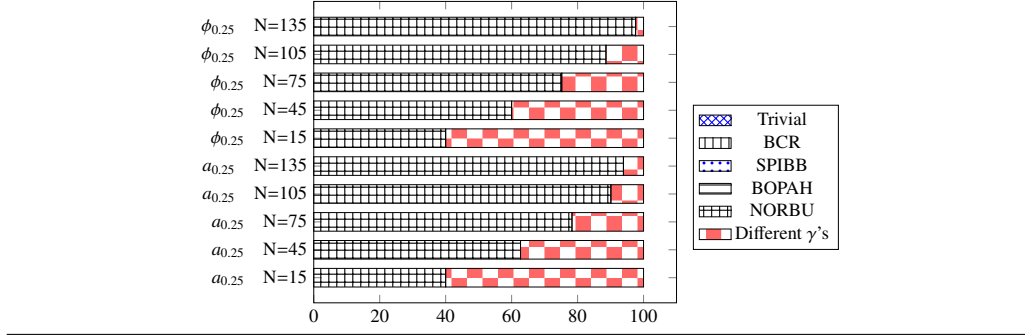
**Figure 6** Policy selection rate by EvC $a_{0.25}$ and EvC $\phi_{0.25}$ in RFL (aggregate of Env. 1, 2, 3 and 4) for different batch sizes.



increases (Ring and Chain) or 2) NORBU (in the RFL environment). This is reasonable since model uncertainty decreases with $N$ and the trivial model will be closer and closer to the true one. We suppose that for not-so-small environments (RFL) the trivial policy can not be trusted for small batch sizes while NORBU manages to cut the posterior space in ambiguity sets that are efficiently optimized over. The policies computed through SPIBB and BOPAH are never selected. Remember that those are stochastic policies that were obtained by improving the batch collector policy that was uniformly random over the actions. Stochastic policies seem not to provide good risk-aware estimates with respect to risk-aware BMDP criteria defined in Eq. (11) and also require sampling more models in order, for the method, to estimate a quantile with the needed accuracy.

Another interesting effect reported in Ring is that for $N = 8$ the trivial policy is picked a considerable amount of times. Both in Ring and in Chain EvC selects more often the output of BCR rather than that of NORBU, even though NORBU in the end is slightly the most performing according to Table 2. In RFL only the policies computed through solving different models sampled from the posterior with different $\gamma$'s and NORBU are selected. The first kind of policy is preferred when the batch is very small $N = 15$, however, the ratio inverts already for $N = 45$ with NORBU that gets more and more chosen with $N$ growing. Both the Trivial policy and the one returned by BCR are always discarded, stressing the superiority of NORBU in this environment typology. Surprisingly, EvC never selects SPIBB nor BOPAH not even in RFL despite its good performance. This is due probably to the difficulty in estimating the quantiles of the performance of a non-deterministic policy such as the output of SPIBB. The algorithm would require a number of sampled models higher than the bail-out hyperparameter.

## 6. Conclusion and future work

This work presents EvC, a method to first evaluate and then select the best risk-aware policies within a set of candidate policies in the context of Offline solutions to Risk-aware Bayesian MDPs. The Risk-aware BMDP defines an elegant mathematical framework that balances the exploitation-caution trade-off in offline model-based sequential decision-making under uncertainty. The set of candidate policies exploited by EvC contains the strategies obtained by solving not only the trivially learned MDP but also other MDPs with transition dynamics sampled from the Bayesian posterior (e.g. the one shown in Eq. (4)) using different discount factors and the solutions of current offline MDP and RL solvers (SPIBB, BOPAH, BCR, NORBU). The estimate

of risk in the presented algorithm provides a probabilistic guarantee for the actual performance of the resulting policy described in Theorem 1 and Theorem 2. The selected solution maximizes the risk-aware utility function of Eq. (11). Since EvC is based on the parallel resolution of a great number of models sampled from the Bayesian posterior we doubt that it could efficiently scale to select policies for MDPs with a great number of states and actions. However, the presented approach should be considered a valuable tool to be exploited for real-world problem-solving through MDP modeling. In such a case time is an affordable resource since the safety of possible humans in the loop would be the priority.

In the future, we aim to improve EvC's method of generation of the set of candidate policies. An interesting direction consists in incrementally enriching the set of candidate policies following some kind of heuristics, e.g. policy improvement by genetic algorithms. An extension to compute robust policies for data-driven POMDPs could be envisaged whether a consistent representation of the model uncertainty can be formalized.

## Acknowledgments

## Code availability

The code for the experiments is open and available in the Github repository: https://github.com/giorgioangel/evc

## References

[1] A. Jonsson, Deep Reinforcement Learning in Medicine, Kidney Diseases 5 (2018) 1–5. `doi:10.1159/000492670`.

[2] B. Mirchevska, C. Pek, M. Werling, M. Althoff, J. Boedecker, High-level Decision Making for Safe and Reasonable Autonomous Lane Changing using Reinforcement Learning, in: 2018 21st International Conference on Intelligent Transportation Systems (ITSC), 2018, pp. 2156–2162. `doi:10.1109/ITSC.2018.8569448`.

[3] R. Laroche, P. Trichelair, R. T. Des Combes, Safe Policy Improvement with Baseline Bootstrapping, in: International Conference on Machine Learning, PMLR, 2019, pp. 3652–3661.

[4] S. Fujimoto, E. Conti, M. Ghavamzadeh, J. Pineau, Benchmarking Batch Deep Reinforcement Learning Algorithms (2019). `arXiv:1910.01708`.

[5] A. Kumar, J. Fu, M. Soh, G. Tucker, S. Levine, Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction, in: Advances in Neural Information Processing Systems 32, 2019, pp. 11784–11794.

[6] Y. Wu, G. Tucker, O. Nachum, Behavior Regularized Offline Reinforcement Learning (2019). `arXiv:1911.11361`.

[7] B. Lee, J. Lee, P. Vrancx, D. Kim, K.-E. Kim, Batch Reinforcement Learning with Hyperparameter Gradients, in: International Conference on Machine Learning, PMLR, 2020, pp. 5725–5735.

[8] J. Chen, N. Jiang, Information-Theoretic Considerations in Batch Reinforcement Learning, in: Proceedings of Machine Learning Research, Vol. 97, PMLR, Long Beach, California, USA, 2019, pp. 1042–1051.

[9] T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, T. Ma, MOPO: Model-based Offline Policy Optimization, in: Advances in Neural Information Processing Systems, Vol. 33, 2020, pp. 14129–14142.

[10] R. Kidambi, A. Rajeswaran, P. Netrapalli, T. Joachims, MOReL: Model-Based Offline Reinforcement Learning, in: Advances in Neural Information Processing Systems, Vol. 33, 2020, pp. 21810–21823.

[11] S. Levine, A. Kumar, G. Tucker, J. Fu, Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems (2020). `arXiv:2005.01643`.

[12] A. Kumar, A. Zhou, G. Tucker, S. Levine, Conservative q-learning for offline reinforcement learning, Advances in Neural Information Processing Systems 33 (2020) 1179–1191.

[13] T. L. Paine, C. Paduraru, A. Michi, C. Gulcehre, K. Zolna, A. Novikov, Z. Wang, N. de Freitas, Hyperparameter selection for offline reinforcement learning, arXiv preprint arXiv:2007.09055 (2020).

[14] R. Munos, C. Szepesvári, Finite-time bounds for fitted value iteration., Journal of Machine Learning Research 9 (5) (2008).

[15] H. Le, C. Voloshin, Y. Yue, Batch policy learning under constraints, in: International Conference on Machine Learning, PMLR, 2019, pp. 3703–3712.

[16] S. Zhang, N. Jiang, Towards hyperparameter-free policy selection for offline reinforcement learning, Advances in Neural Information Processing Systems 34 (2021).

[17] C.-H. H. Yang, Z. Qi, Y. Cui, P.-Y. Chen, Pessimistic model selection for offline deep reinforcement learning, arXiv preprint arXiv:2111.14346 (2021).

[18] A. Nilim, L. El Ghaoui, Robust Control of Markov Decision Processes with Uncertain Transition Matrices, Operations Research 53 (5) (2005) 780–798. arXiv:https://doi.org/10.1287/opre.1050.0216, doi:10.1287/opre.1050.0216.

[19] G. N. Iyengar, Robust Dynamic Programming, Mathematics of Operations Research 30 (2) (2005) 257–280. arXiv:https://doi.org/10.1287/moor.1040.0129, doi:10.1287/moor.1040.0129.

[20] E. Delage, S. Mannor, Percentile optimization for markov decision processes with parameter uncertainty, Operations research 58 (1) (2010) 203–213.

[21] M. Petrik, M. Ghavamzadeh, Y. Chow, Safe Policy Improvement by Minimizing Robust Baseline Regret, Advances in Neural Information Processing Systems 29 (2016) 2298–2306.

[22] M. Petrik, R. H. Russel, Beyond Confidence Regions: Tight Bayesian Ambiguity Sets for Robust MDPs, in: Advances in Neural Information Processing Systems 32, Vol. 32, 2019.

[23] B. Behzadian, R. Hasan Russel, M. Petrik, C. Pang Ho, Optimizing Percentile Criterion using Robust MDPs, in: Proceedings of The 24th International Conference on Artificial Intelligence and Statistics, Vol. 130 of Proceedings of Machine Learning Research, PMLR, 2021, pp. 1009–1017.

[24] E. A. Lobo, M. Ghavamzadeh, M. Petrik, Soft-robust algorithms for batch reinforcement learning (2021). arXiv:2011.14495.

[25] N. Jiang, A. Kulesza, S. Singh, R. Lewis, The dependence of effective planning horizon on model accuracy, in: Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems, 2015, pp. 1181–1189.

[26] D. Precup, Eligibility traces for off-policy policy evaluation, Computer Science Department Faculty Publication Series (2000) 80.

[27] Y. Chandak, S. Niekum, B. C. da Silva, E. Learned-Miller, E. Brunskill, P. S. Thomas, Universal off-policy evaluation, Advances in Neural Information Processing Systems 34 (2021).

[28] M. Strens, A Bayesian framework for Reinforcement Learning, in: In Proceedings of the Seventeenth International Conference on Machine Learning, ICML, 2000, pp. 943–950.

[29] A. Sharma, J. Harrison, M. Tsao, M. Pavone, Robust and Adaptive Planning under Model Uncertainty, Proceedings of the International Conference on Automated Planning and Scheduling 29 (1) (2019) 410–418.

[30] S. Depeweg, J. M. Hernández-Lobato, F. Doshi-Velez, S. Udluft, Uncertainty decomposition in bayesian neural networks with latent variables, arXiv preprint arXiv:1706.08495 (2017).

[31] S. Depeweg, J.-M. Hernandez-Lobato, F. Doshi-Velez, S. Udluft, Decomposition of uncertainty in bayesian deep learning for efficient and risk-sensitive learning, in: International Conference on Machine Learning, PMLR, 2018, pp. 1184–1193.

[32] L. N. Steimle, D. L. Kaufman, B. T. Denton, Multi-model markov decision processes, IISE Transactions 53 (10) (2021) 1124–1139.

[33] J. R. Rice, The algorithm selection problem, in: Advances in computers, Vol. 15, Elsevier, 1976, pp. 65–118.

[34] R. Laroche, R. Feraud, Reinforcement learning algorithm selection, in: International Conference on Learning Representations 2018, 2018.

[35] A. Majumdar, M. Pavone, G. Hager, S. Thomas, M. Torres-Torriti, How Should a Robot Assess Risk? Towards an Axiomatic Theory of Risk in Robotics, in: Robotics Research, Springer International Publishing, Cham, 2020, pp. 75–84.

[36] Mausam, A. Kolobov, Planning with Markov Decision Processes: An AI perspective, Synthesis Lectures on Artificial Intelligence and Machine Learning 6 (1) (2012) 1–210.

[37] P. Artzner, F. Delbaen, E. Jean-Marc, D. Heath, Coherent Measures of Risk, Mathematical Finance 9 (1999) 203 – 228. doi:10.1111/1467-9965.00068.

[38] R. T. Rockafellar, S. Uryasev, Conditional value-at-risk for general loss distributions, Journal of banking & finance 26 (7) (2002) 1443–1471.

[39] K. Briggs, F. Ying, How to estimate quantiles easily and reliably, Mathematics Today 2018 (February) (2018) 26–29.

[40] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, W. Zaremba, OpenAI Gym (2016). arXiv:1606.01540.