
Conservative Value Priors: A Bayesian Path to Offline Reinforcement Learning

Filippo Valdettaro¹ Yingzhen Li¹ A. Aldo Faisal^{1,2,3}

¹Department of Computing, Imperial College London, UK

²Department of Bioengineering, Imperial College London, UK

³Chair in Digital Health & Data Science, University of Bayreuth, Germany

Abstract

Offline reinforcement learning (RL) seeks to improve a policy using only a fixed dataset of past interactions, without further environment exploration. To avoid overly optimistic decisions on uncertain or out-of-distribution actions, the learned policy should be supported by the data. Existing methods typically enforce this via heuristic modifications to objectives or value functions that encourage more conservative action selection. In contrast, we propose a principled alternative by introducing a conservative value prior, thereby modelling the belief that policies are expected to perform poorly unless the behavioural data provides evidence to the contrary. This yields a posterior that assigns high value only to supported actions, guiding the agent toward policies grounded in the data. Our approach thereby unifies Bayesian decision-making, uncertainty quantification and value regularisation while effectively mitigating distributional shift in offline RL. We develop this framework in a model-free setting for continuous control in deterministic environments. We first present an exact inference algorithm for small-scale problems, then extend it to a scalable deep learning variant compatible with standard off-policy algorithms. Our method achieves strong performance on benchmark locomotion tasks, outperforming comparable model-free baselines thanks to the milder yet effective form of regularisation employed.

1 Introduction

Offline reinforcement learning (RL) can harness datasets of suboptimal demonstrations to learn effective policies without direct interaction or additional exploration on an environment. This extends the potential for the application of RL in domains where simulation is not realistically possible, where a large amount of operational data have already been collected and where direct suboptimal interaction with the underlying environment can be dangerous or costly, such as healthcare [Gottesman et al., 2019, Komorowski et al., 2018], robotics [Sinha et al., 2022, Kalashnikov et al., 2018, Dasari et al., 2020, Kendall et al., 2019] or recommender systems [Huang et al., 2022, Xiao and Wang, 2021]. However, a lack of interaction with the environment comes with additional challenges. In particular, traditional off-policy actor-critic algorithms can perform very poorly when naively applied to an offline dataset: without appropriate regularisation, uncontrolled extrapolation can cause the critic to assign high values to actions that it has never encountered in the dataset, leading to the actor selecting policies that are entirely unsupported by the data.

A flexible strategy to mitigate this issue is to train a critic that learns a *conservative* or *pessimistic* value function by modifying the value function estimate: one approach is to incorporate a penalty on state-action pairs with high uncertainty into the value estimate, often measured by the disagreement among an ensemble of critics [Ghasemipour et al., 2022, An et al., 2021]. While ensemble methods can be effective, they tend to be computationally intensive, and the extent of the uncertainty penalty

may need to be determined in an *ad hoc* manner. Other approaches have brought about methods that bypass explicit uncertainty quantification by training the critic to estimate a lower bound on the value function based solely on the observed data, aiming to provide robustness against worst-case outcomes [Kumar et al., 2020, Cheng et al., 2022]. However, such worst-case formulations may be excessively pessimistic, potentially hindering overall performance [Xu and Mannor, 2006].

Here, we explore an alternative approach to learning a conservative value function estimation for offline RL rooted in Bayesian inference. We build on the proof-of-concept work in Valdetaro and Faisal [2024], extend it to continuous control, and provide a scalable deep learning variant of the method. Our approach is inspired by Bayesian decision theory: the actor is trained to optimise a posterior expectation of utility (value), which is a provably optimal and principled decision criterion [Robert et al., 2007], and is in general preferable to worst-case robustness as the latter can significantly harm average performance [Xu and Mannor, 2006]. Choosing a prior belief over value functions that is conservative will naturally induce conservatism in the posterior value function estimate of actions outside the dataset’s support. This can be interpreted as a less extreme form regularisation than in algorithms that learn worst-case lower bounds of value functions. Thus, we train a critic with a conservative value prior (CVP), so that after inference only those state-actions that are supported by the data will be assigned high values. We achieve this in an off-policy model-free way in deterministic environments by structuring inference so that it respects the temporal-difference (TD) structure of a Markovian environment to get consistent value functions after conditioning on the observed data. As a natural by-product of our inference scheme, we obtain posterior epistemic uncertainty estimates, although we keep these decoupled from decision-making.

Our main contributions are threefold and summarised as follows: 1. We develop the CVP probabilistic formulation for continuous control in deterministic environments with exact inference showcasing qualitative results and uncertainty visualisations on a classical control environment in the low-data regime. 2. We then propose a scalable, deep learning-based variant of CVP for offline RL, achieved by adding an easily-computable term to the critic’s loss, which constitutes a simple modification applicable to any off-policy algorithm. 3. We evaluate our method on the D4RL [Fu et al., 2020] robotics locomotion tasks, on which we find it performs better than previous comparable methods and recovers similar performance to more algorithmically complex methods.

2 Related work

In the following, we introduce and motivate the elements of conservative value estimation, existing approaches, and model-free inference in RL. Additionally, we cover function-space inference, which enables Bayesian reasoning over function outputs rather than parameters, which is necessary to introduce a conservative prior over value.

Conservative value estimation. A core challenge in offline RL is avoiding distributional shift between the learned and behaviour policies. One approach adds constraints to keep the actor close to the behaviour policy, either explicitly [Fujimoto and Gu, 2021] or implicitly [Nair et al., 2020], but this is sensitive to suboptimal data. Instead, faithful generative models of the behaviour policy can constrain actions to merely by in the dataset’s support [Fujimoto et al., 2019, Zhou et al., 2021], though they rely on accurately modelling the behaviour policy, which is itself challenging. Unlike behaviour regularisation, conservative value estimation methods are less sensitive to poor-quality data. Uncertainty-based approaches [An et al., 2021, Ghasemipour et al., 2022, Yang et al., 2024] use pessimistic targets from ensembles to prevent value overestimation in uncertain regions. However, ensembles are computationally expensive and lack strong theoretical guarantees despite empirical success [Lakshminarayanan et al., 2017, D’Angelo and Fortuin, 2021]. Beyond uncertainty-based methods, critic regularisation can also use action discrepancy penalties [Wu et al., 2019, Tarasov et al., 2024, Kostrikov et al., 2021], or directly learn conservative value lower bounds robust to worst-case scenarios [Kumar et al., 2020, Lyu et al., 2022, Cheng et al., 2022]. Our work builds on Valdetaro and Faisal [2024], where conservatism arises from Bayesian inference with a conservative value prior, ensuring high values are assigned only to supported regions, thus naturally avoiding OOD actions without explicit robustness.

Bayesian Inference in RL. Bayesian methods offer an approach to the exploration-exploitation challenge of online RL, either through uncertainty-guided heuristics [O’Donoghue et al., 2018,

Agrawal and Goyal, 2017] or formalised as expected value maximisation in belief-space MDPs [Duff, 2002, Poupart et al., 2006]. While prior uncertainty-based offline RL methods emphasise uncertainty estimation and conservative value learning, a principled posterior expected value maximisation formulation of offline RL remains an open challenge. Our work proposes a step towards such a formulation for offline, model-free RL that retains the key Bayesian principle of posterior value optimisation.

A central challenge in off-policy model-free Bayesian inference setting is that value targets are not observed directly. Unlike supervised learning, or even on-policy Bayesian RL (e.g., GP-SARSA [Engel et al., 2005]), the agent must infer values by leveraging the Markov structure and stitching together information from transitions. While ensemble-based methods [Lakshminarayanan et al., 2017, Ovadia et al., 2019, Osband et al., 2018] have been used to approximate epistemic uncertainty in RL, their theoretical grounding and uncertainty quantification capabilities specifically in the context of TD learning remains unsettled [Osband et al., 2018, Touati et al., 2020]. In contrast, we adopt a fully probabilistic approach, placing a prior over value functions and performing inference consistent with the Bellman equation and the data.

Inference in function space Our method relies on placing a prior directly on the value function. Thus, we need to carry out Bayesian inference in function space rather than parameter space, which comes with an additional set of complexities. While some nonparametric approaches, such as Gaussian processes (GPs), naturally carry out inference in function space, these require approximations to be scalable to large datasets [Hensman et al., 2013, Titsias, 2009] and classically rely on hand-crafted kernels which is undesirable for complex tasks or environments. The field of functional variational inference [Burt et al., 2020, Ma and Hernández-Lobato, 2021] seeks to approximate inference in value-function space while still employing expressive parametric models. The deep learning variant of our approach is closest to the framework for functional variational inference proposed in Hafner et al. [2020], where the prior knowledge is injected into training by sampling pseudo-data points from a prior distribution.

3 Conservative Value Priors

Setting and notation. We work with the discounted Markov Decision Process (MDP) formalism [Sutton et al., 1998] with states $s \in \mathcal{S}$, actions $a \in \mathcal{A}$, discount factor γ , rewards $r \in \mathbb{R}$ sampled from an unknown reward function with mean $R(s, a)$ and next states $s' \in \mathcal{S}$ sampled from an unknown transition kernel $P(s'|s, a)$. For tractability, we assume deterministic transitions in the analysis. The objective is to use a fixed dataset of transitions $(s_i, a_i, r_i, s'_i) \in \mathcal{D}$ to learn a policy $\pi(a|s)$ mapping states to distributions over actions (with shorthand $a = \pi(s)$ for deterministic policies) that maximises expected cumulative discounted returns $\mathbb{E}(\sum_{t=0}^{\infty} \gamma^t r_t)$ where r_t is the reward sampled at step t , with $s_{t+1} \sim P(\cdot|s_t, a_t)$, $a_{t+1} \sim \pi(s_t)$, $s_0 \sim \eta_0(s)$ with η_0 being some unknown initial state distribution. We refer to the (discounted) state distribution under π as $\eta^\pi(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi)$. Finally, for a given policy π we denote $V^\pi(s) = \mathbb{E}(\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s)$ and $Q^\pi(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}(V^\pi(s'))$.

Conservative Value Priors. Here, we define and give motivation for conservative value priors.

Definition 1. A prior over state-action values $p(Q^\pi(s, a))$ is defined as conservative with respect to policy π if $\mathbb{E}(Q^\pi(s, a)) < V^\pi(s)$ for any $a \in \mathcal{A}$ and all $s \in \text{supp}(\eta^\pi)$, where the expectation is taken with respect to the prior p .

The policy π is the relevant policy that we wish to learn, which in the offline RL case can be taken to be the best policy in the support of the dataset. This property is crucial in a Bayesian treatment of offline RL where datasets have limited cover. Suppose in some state s the dataset provides no information about action $a \neq \pi(s)$ but does allow an accurate estimation of $V^\pi(s)$. Then, the posterior belief about $Q^\pi(s, a)$ will remain similar to the prior belief, so that $\mathbb{E}(Q^\pi(s, a) | \mathcal{D}) \approx \mathbb{E}(Q^\pi(s, a))$. If this prior is not conservative by Definition 1, since $\mathbb{E}(Q^\pi(s, a)) > V^\pi(s) = Q^\pi(s, \pi(s))$, a Bayesian posterior value maximisation framework will lead to choosing the OOD $a = \text{argmax}_{a'} \mathbb{E}(Q^\pi(s, a') | \mathcal{D})$, over the supported $\pi(s)$. Therefore, a conservative prior must be in place to avoid mistakenly favouring unsupported actions. The condition that s is in the support of π ensures this property only needs to hold for states that are relevant towards decision-making. We visualise this argument in Fig. 1.

Choosing a suitable prior. In practice, we will need to choose prior that satisfies Definition 1. We propose two approaches for this:

1. **Manually choose mean.** Often, we can leverage known properties of the specific RL set up to define a prior with suitable mean. While this requires task-specific knowledge, many tasks have a natural baseline value that can be used as a baseline conservative prior mean. See, for example, the robotics experiments in Section 5.2.
2. **Reward preprocessing (RP).** Translate all rewards in the dataset by a constant $r \rightarrow r - r_{\min}$, with r_{\min} being the smallest observed reward in the dataset. A value of 0 then corresponds to observing the worst-case reward outcome at every step, thus likely making a zero-mean prior conservative for typical datasets and environments.

While RP provides a general recipe for finding a suitable prior mean, the drawback is that an excessively pessimistic prior may introduce excessive regularisation and harm performance, so the choice of approach should be considered on a task-by-task basis.

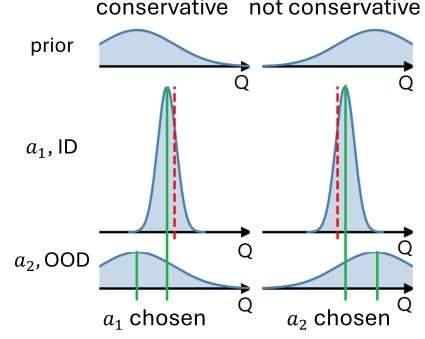


Figure 1: Effect of conservative (left) vs non-conservative (right) priors on action selection for an in-distribution (ID) action a_1 and an OOD action a_2 . The available data (red line) leaves the belief of OOD actions unchanged, so if the prior is not conservative by Definition 1, maximising posterior expected value (green line) will lead the agent to choose the OOD action a_2 .

4 CVP via Gaussian Processes

We present here an algorithm that implements the CVP idea with GPs with exact inference, suitable for simple environments and small datasets. Through a qualitative evaluation of its decision-making and uncertainty quantification capabilities in the mountain car classic control task, we demonstrate the soundness of the CVP approach in its most direct form.

4.1 Off-policy evaluation with GPs

The derivation we present here builds on the work in Valdetaro and Faisal [2024], which itself provides an off-policy generalisation of the line of work in Engel et al. [2005], to achieve TD off-policy inference in value-function space for continuous state and action spaces. Our objective is to find a policy that maximises the off-policy objective [Degris et al., 2012]

$$J(\pi) = \int \rho(s) V^\pi(s) ds \approx \frac{1}{|\mathcal{D}|} \sum_{s \in \mathcal{D}} Q^\pi(s, \pi(s)), \quad (1)$$

where we have approximated the state distribution under the dataset policy ρ with the dataset’s empirical state distribution. Since Q^π is an unknown quantity, J is also unknown, and Bayesian decision theory [Robert et al., 2007] dictates that we maximise its expected posterior value, so the task is now to find an expression for $\mathbb{E}(Q^\pi | \mathcal{D})$, which we can then use to find a π that optimises Eq.1.

A key conceptual challenge in Bayesian model-free off-policy evaluation is that values are not observed directly. They must instead be inferred from information obtained in individual transitions and the knowledge that the environment is Markovian. By including the Bellman equation directly into the inference process, we ensure that posterior samples of the value function respect the environment’s Markov structure and consistently stitch information about value across multiple transitions. For a deterministic policy π and a deterministic environment, the action-values at any s, a satisfy the Bellman equation

$$Q^\pi(s, a) = R(s, a) + \gamma Q^\pi(s', \pi(s')), \quad (2)$$

where s' is the state deterministically reached from state s after taking action a . Thus, the probabilistic model we employ assumes that the observed (reward samples r) and latent (action-values Q) variables

relevant to the transitions in the dataset $(s_i, a_i, r_i, s'_i) \in \mathcal{D}$ are related through

$$r_i = Q^\pi(s_i, a_i) - \gamma_i Q^\pi(s'_i, \pi(s'_i)) + \varepsilon_i, \quad (3)$$

where $\gamma_i = \gamma$ if s'_i is not terminal and 0 otherwise, and ε_i a small zero-mean independent Gaussian noise term with variance σ_r^2 . Next, we consider a given (conservative) Gaussian prior on Q with constant prior mean μ_Q and covariance kernel that factorises across state and action spaces $k_Q((s_1, a_1), (s_2, a_2)) = k_s(s_1, s_2)k_a(a_1, a_2)$.

Since each r_i is itself a linear combination of Gaussian random variables, we can establish its prior mean and the covariance between any two distinct observed rewards:

$$\mathbb{E}(r_i) = (1 - \gamma_i)\mu_Q \quad (4)$$

$$\text{cov}(r_i, r_j) = k(x_i, x_j) - \gamma_j k(x_i, x'_j) - \gamma_i k(x'_i, x_j) + \gamma_i \gamma_j k(x'_i, x'_j), \quad (5)$$

where we used the shorthand notation $x_k = (s_k, a_k)$ and $x'_k = (s'_k, \pi(s'_k))$. We can also compute the prior covariance between the observed rewards and the action-value of any arbitrary state-action:

$$\text{cov}(Q^\pi(s, a), r_i) = k((s, a), x_i) - \gamma_i k((s, a), x'_i). \quad (6)$$

We can now find the posterior value distribution for any state-action s, a by using the Gaussian conditioning formulas for posterior mean μ^* and covariance Σ^* [Rasmussen et al., 2006]:

$$\mu^* = \mathbf{m}_r + \mathbf{K}_{QR}^\top (\mathbf{K}_R + \sigma_r^2 \mathbf{I})^{-1} (\mathbf{r} - \mathbf{m}_r) \quad (7)$$

$$\Sigma^* = \mathbf{K}_Q - \mathbf{K}_{QR}^\top (\mathbf{K}_R + \sigma_r^2 \mathbf{I})^{-1} \mathbf{K}_{QR}, \quad (8)$$

where we can populate \mathbf{m}_r , \mathbf{K}_R , \mathbf{K}_{QR} and \mathbf{K}_Q with the entries $\mathbb{E}(r_i)$, $\text{cov}(r_i, r_j)$, $\text{cov}(Q^\pi(s, a), r_i)$ and $k((s, a), (s, a))$ respectively. Notice that the action-value posterior’s dependence on policy is differentiable and present through the terms containing $x' = (s', \pi(s'))$.

This allows us to compute posterior mean (and variance) for any $Q^\pi(s, a)$, giving a differentiable expression for the posterior expectation of Eq. 1 in terms of π , which then be optimise through gradient ascent. This is summarised in Alg. 1. In principle, the posterior variance could also be used for action selection, but keeping in line with Bayesian decision theory we only consider posterior mean as the objective to maximise.

Algorithm 1 CVP for offline RL with GPs

Require: Dataset \mathcal{D} , conservative Gaussian prior on Q^π , parametric actor $\pi_\theta : \mathcal{S} \rightarrow \mathcal{A}$, learning rate η

while not converged **do**

$\hat{Q} \leftarrow \mathbb{E}(Q^{\pi_\theta} | \mathcal{D})$ ▷ Evaluate $\mathbb{E}(Q^{\pi_\theta} | \mathcal{D})$, Eq. 7

$J(\theta) \leftarrow \frac{1}{|\mathcal{D}|} \sum_{s \in \mathcal{D}} \hat{Q}(s, \pi_\theta(s))$ ▷ Posterior mean of objective in Eq. 1

$\theta \leftarrow \theta + \eta \nabla_{\theta'} J(\theta')$ ▷ Change π_θ in direction of increasing J .

end while

4.2 Continuous control with small datasets

We qualitatively evaluate the exact inference CVP approach on the classic mountain car control problem [Moore, 1990], demonstrating its ability to learn good offline policies and meaningful uncertainty estimates. This environment is chosen for its two-dimensional state space, which allows intuitive visualisation of posterior value and uncertainty. The task is set in a low-dimensional state-action space and small-data regime. A brief description of the environment is provided in Appendix A. We use a variant with continuous actions and a sparse reward: 1 if the agent reaches the goal, and 0 otherwise. This experiment showcases three core capabilities of our method. First, it is able to stitch information from different transitions in a dataset to obtain policies that lead to in-distribution trajectories with high rewards. Secondly, it correctly assigns low values to policies that rely on OOD state-actions. Thirdly, it produces consistent Bayesian uncertainty estimates over multiple steps.

We choose the prior mean to be $\mu_Q = 0$, following the RP guideline from Section 3. Since all rewards are non-negative for this environment and the goal is reachable from any initial state, it follows that

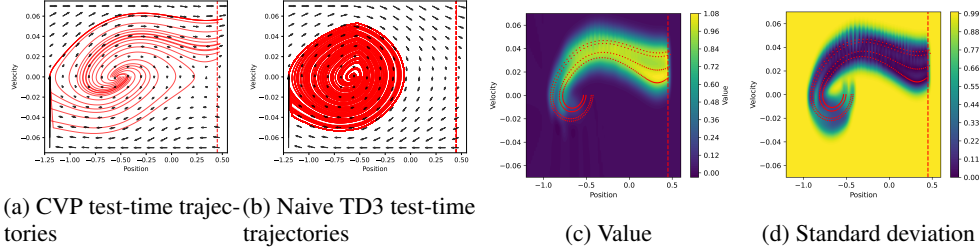


Figure 2: Mountain car state-space visualisation. The black arrows in Figs. 2a and 2b show the transitions in state-space after following the learned policies for one timestep and we also include 15 trajectories that arise from uniformly spaced starting states (red continuous lines). Figs. 2c and 2d display learned posterior values and uncertainty (standard deviation) of the trained CVP agent with exact inference. The red dots indicate the states observed in the dataset. In all figures, crossing the vertical red dotted line corresponds to reaching the goal.

$0 < V^\pi(s) < 1$ in the support of η^π and therefore $\mu_Q = 0$ satisfies Definition 1. The RBF covariance function’s lengthscale was selected manually (see Appendix A for further details).

We apply Alg. 1 to a dataset consisting of 4 expert trajectories from a pre-trained agent (shown as red dots in Fig. 2). The resulting values, uncertainty and trajectories from the learned CVP policy are shown in Fig. 2. In Fig. 2a we display the dynamics under this policy and full episodic trajectories from uniformly spaced initial states, all of which successfully reach the goal. In contrast, a naive TD3 agent trained on the dataset as the replay buffer fails to do so (Fig. 2b). Further experimental details are given in Appendix A.

This behaviour illustrates the desired properties outlined at the start Section 4.2. The posterior value is correctly propagated across steps, being approximately equal to 1 close to the goal and slowly decaying as more steps become necessary to reach it. High values are restricted to in-distribution regions where data supports successful goal-reaching behaviour; elsewhere, values remain low. Interestingly, even in OOD regions where the policy successfully generalises to solve the task, it still assigns a low posterior value if it is out of the data distribution: the agent ‘tries its best’ but does not expect a high value unless the data supports it. In a complementary fashion, posterior uncertainty is low in the in-distribution regions and high OOD, as expected by a consistent quantification of epistemic uncertainty.

5 CVPs via Neural Networks

While exact Bayesian inference is suitable for small datasets and simple environments where a hand-crafted kernel is easily specified, sequential decision making from larger datasets on complex environments requires a more scalable approach. One option is to approximate exact inference with variational inference [Titsias, 2009], but this still leaves the difficult task of finding adequate prior kernels, and methods that try to learn these [Wilson et al., 2016, Ober et al., 2021, van Amersfoort et al., 2021] require introducing a highly non-trivial layer of algorithmic complexity on top of the already challenging task of offline RL. Therefore, here we opt instead to forego approximations to the full posterior distribution of a GP altogether and instead focus on approximately regressing the posterior mean only, as parametrised by a neural network. To implement this, we take inspiration from the approach in Hafner et al. [2020], where the inductive bias of inference in function space with a prior is introduced by sampling pseudo-data points from the prior.

5.1 A scalable implementation of CVPs

The key quantity that is optimised in Eq. 1 for decision making in Section 4 is the posterior mean of Q^π . To maintain scalability to large datasets and complex environments, we employ a neural network $\hat{Q}_\theta(s, a)$ to represent it. For tractability, we simplify the probabilistic model from Section 4 to employ independent Gaussian priors across state-actions, corresponding to choosing a delta covariance kernel $k_Q((s_1, a_1), (s_2, a_2)) = \delta_s(s_1 - s_2)\delta_a(a_1 - a_2)$. We can then find a training objective that ensures that \hat{Q} regresses to this posterior mean. We show in Appendix B that the appropriate loss function to

achieve this is given by

$$L(\theta) = L_B(\theta) + \alpha \mathbb{E}_{(s,a) \sim p_{\text{pseudo}}(s,a)} (\hat{Q}_\theta(s,a) - \mu_Q)^2, \quad (9)$$

where L_B is the standard Bellman TD critic loss

$$L_B(\theta) = \sum_{(s,a,r,s') \in \mathcal{D}} (\hat{Q}_\theta(s,a) - r - \gamma \hat{Q}_\theta(s', \pi(s')))^2, \quad (10)$$

μ_Q is the prior mean on Q , p_{pseudo} is a prior pseudo-data distribution used to regularise the learned value function and α is a hyperparameter. The result in Appendix B holds when p_{pseudo} is uniform, where we show the L_B term can be interpreted as a log-likelihood loss and the other term as the contribution from the prior, with the corollary that the regularisation strength α can be interpreted as being inversely proportional to the prior variance. An intuitive interpretation of Eq. 9 is that \hat{Q} will be incentivised to fit to the data in the in-distribution regions where the L_B term dominates and will instead regress to the prior otherwise. We further note that the Eq. 9 can be interpreted as a TD analogue of Eq. 2 in Hafner et al. [2020], where the KL between two Gaussians recovers the same prior term proportional to the difference of the means squared.

While the assumption of independent state-actions would not lead to useful generalisation with non-parametric methods in continuous state-action spaces (such as GPs), when employing a parametric function approximator with strong generalisation capabilities, such as neural networks, this is not problematic: since hand-crafted kernels are unlikely to achieve good results in complicated environments anyway, we let the inductive bias implicit in the neural networks handle generalisation instead. On the other hand, the condition that p_{pseudo} is taken to be uniform will run into serious issues in high-dimensional environments. Thus, we instead restrict ourselves to regularising the most relevant state-actions towards decision-making. For a dataset state distribution ρ , a natural pseudo-data distribution for a stochastic policy would be $\rho(s)\pi(a|s)$ whereas for a deterministic policy we can sample so as to ensure regularisation around the learned policy’s boundary

$$p_{\text{pseudo}}(s,a) = \rho(s)\mathcal{N}(a|\pi_\theta(s), \sigma_a^2), \quad (11)$$

for some constant σ_a^2 with similar intuition that this prioritises regularising the most relevant actions, which are those closest to the ones considered by the learned policy, as in Hafner et al. [2020]. The resulting critic loss is summarised in Algorithm 2, which we observe is a straightforward modification of standard off-policy online algorithms’ loss that simply involves introducing a readily-computable regularisation term arising from the prior’s contribution.

Algorithm 2 Scalable CVP critic update loss

Require: Batch of transitions \mathcal{B} sampled from a dataset \mathcal{D} , online off-policy algorithm with deterministic or stochastic actor π , critic Q_θ and critic loss $L_B(\theta)$, conservative prior mean μ_Q , prior regularisation strength α , pseudo-data number of samples n and noise scale σ_a .

for $(s, a, r, s') \in \mathcal{B}$ **do**
 for $s_p \in \{s, s'\}$ **do**

Sample n actions $a_p^{(s_p)} \leftarrow \begin{cases} a_p^{(s_p)} = \pi(s_p) + \sigma_a z, z \sim \mathcal{N}(\cdot|0, \mathbf{I}_n) & \text{if } \pi \text{ deterministic} \\ a_p^{(s_p)} \sim \pi(\cdot|s_p) & \text{if } \pi \text{ stochastic} \end{cases}$

end for

end for

$\mathcal{S}_\mathcal{B} \leftarrow \{s, s' \mid (s, a, r, s') \in \mathcal{B}\}$ \triangleright Use empirical state distribution to approximate $\rho(s)$

$$L \leftarrow L_B(\theta) + \alpha \frac{1}{|\mathcal{B}|n} \sum_{s_p \in \mathcal{S}_\mathcal{B}} \sum_{a_p^{(s_p)}} (Q_\theta(s_p, a_p^{(s_p)}) - \mu_Q)^2$$

return loss L

Having established an approach to include the conservative prior’s effect in the critic’s evaluation, we can adapt any off-policy online algorithm by adding the new term in Eq. 9 to obtain a suitably regularised offline RL critic. In Appendix C we present a comparison between applying Alg. 1 and Alg. 2 for both a stochastic (SAC) and deterministic (TD3) base algorithm on a structured, easily interpretable toy environment, and observe that the key desirable behaviour is consistent across all three implementations. For the remainder of this work, we use the stochastic policy variant, with SAC as the base algorithm for experiments.

5.2 Offline Robotic Locomotion

Table 1: D4RL results on random (-r), medium (-m), medium-replay (-m-r), medium-expert (-m-e) and expert (-e) datasets. Scores are normalised such that a random policy scores 0 and an expert policy scores 100. Baseline results are taken from Lyu et al. [2022]. We bold those scores within one standard deviation of the method with the highest average.

Task	BC	SAC	CQL	TD3+BC	IQL	CVP-SAC (ours)
halfcheetah-r	2.2±0.0	29.7±1.4	17.5±1.5	11.0±1.1	13.1±1.3	32.5±1.5
hopper-r	3.7±0.6	9.9±1.5	7.9±0.4	8.5±0.6	7.9±0.2	27.5±8.0
walker2d-r	1.3±0.1	0.9±0.8	5.1±1.3	1.6±1.7	5.4±1.2	6.9±3.0
halfcheetah-m	43.2±0.6	55.2±27.8	47.0±0.5	48.3±0.3	47.4±0.2	66.0±1.8
hopper-m	54.1±3.8	0.8±0.0	53.0±28.5	59.3±4.2	66.2±5.7	72.2±20.0
walker2d-m	70.9±11.0	-0.3±0.2	73.3±17.7	83.7±2.1	78.3±8.7	84.1±1.0
halfcheetah-m-r	37.6±2.1	0.8±1.0	45.5±0.7	44.6±0.5	44.2±1.2	57.2±0.9
hopper-m-r	16.6±4.8	7.4±0.5	88.7±12.9	60.9±18.8	94.7±8.6	100.5±4.2
walker2d-m-r	20.3±9.8	-0.4±0.3	81.8±2.7	81.8±5.5	73.8±7.1	81.5±11.1
halfcheetah-m-e	44.0±1.6	28.4±19.4	75.6±25.7	90.7±4.3	86.7±5.3	95.4±0.6
hopper-m-e	53.9±4.7	0.7±0.0	105.6±12.9	98.0±9.4	91.5±14.3	103.8±10.1
walker2d-m-e	90.1±13.2	1.9±3.9	107.9±1.6	110.1±0.5	109.6±1.0	109.3±0.9
halfcheetah-e	91.8±1.5	-0.8±1.8	96.3±1.3	96.7±1.1	95.0±0.5	95.1±0.6
hopper-e	107.7±0.7	0.7±0.0	96.5±28.0	107.8±7	109.4±0.5	110.7±1.5
walker2d-e	106.7±0.2	0.7±0.3	108.5±0.5	110.2±0.3	109.9±1.2	110.0±0.2
Average	49.6	9.0	67.3	67.6	68.9	76.9

We demonstrate our framework’s scalability by evaluating the SAC-based CVP variant (Alg. 2) on D4RL locomotion tasks [Fu et al., 2020]. We build our implementation on top of the CORL [Tarasov et al., 2022] CQL code, and do not modify any hyperparameters relating to the architecture or optimisation of the base SAC algorithm, with the only exception being the optimiser used for the automatic entropy tuning parameter loss [Haarnoja et al., 2018b]: we found that ADAM [Kingma and Ba, 2014] updates on this single parameter led to similar performance but presented convergence issues after training stabilised, so we used stochastic gradient descent to optimise this loss instead. Full hyperparameter settings and other experimental details are provided in Appendix D.

Choosing CVP hyperparameters. The two main hyperparameters necessary for CVP are prior mean and α . For the choice of prior mean, the locomotion tasks have a reward of the form *movement reward* minus *control cost*, so a ‘do nothing’ policy would have value approximately equal to 0. Thus, setting $\mu_Q = 0$ provides a suitable conservative prior mean. While we cannot prove that this is a CVP in the strict sense of Definition 1, this is a reasonable base value to start from. We found that for one task (hopper-medium-expert), further tuning of μ_Q helped performance while using the RP heuristic (see Section 3) was instead beneficial for other tasks. Choosing α is less straightforward since, unlike μ_Q , there was no clear default, so we tuned it per task. We briefly discuss in Appendix E how training diagnostics can help guide the choice of α .

D4RL results. In Table 1 we display results on a variety of D4RL v2 locomotion tasks [Fu et al., 2020]. We report the results of evaluating the best-performing hyperparameters after evaluating on the ground-truth environments, reporting averages and sample standard deviations over 5 random seeds. The values for the other methods reported are taken from Lyu et al. [2022]. We compare results with behaviour cloning (BC), SAC [Haarnoja et al., 2018a], CQL [Kumar et al., 2020], TD3+BC [Fujimoto and Gu, 2021] and IQL [Kostrikov et al., 2022]. These, like our method, are straightforward adaptations of base off-policy model-free algorithms. We include in Appendix F additional results comparing to other two more sophisticated model-free algorithms that employ additional elaborate ingredients, MCQ [Lyu et al., 2022] and ATAC [Cheng et al., 2022], where we see that our simpler algorithm still recovers similarly good performance, achieving 97% and 102% of their averaged D4RL score respectively.

Algorithmically, the most closely related baseline to our approach in Table 1 is CQL. CVP-SAC shares the same base algorithm and the difference between the two boils down to the form of the critic regularisation term, a consequence of the two approaches’ fundamentally different starting principles. We see a significant improvement in performance from these differences.

5.3 Hyperparameter study

To illustrate the impact of α and the choice of prior mean on the performance of our algorithm, we present results of varying these on the exemplar halfcheetah-medium task. The overall trend we notice is that, generally speaking, decreasing regularisation strength (decreasing α or increasing μ_Q) tends to have a beneficial effect on performance up to a critical point, after which performance and training becomes unstable. On the other hand, adding excessive regularisation tends to have a more gradual impact on performance. These insights can help tune α from training diagnostics, as discussed in Appendix E.

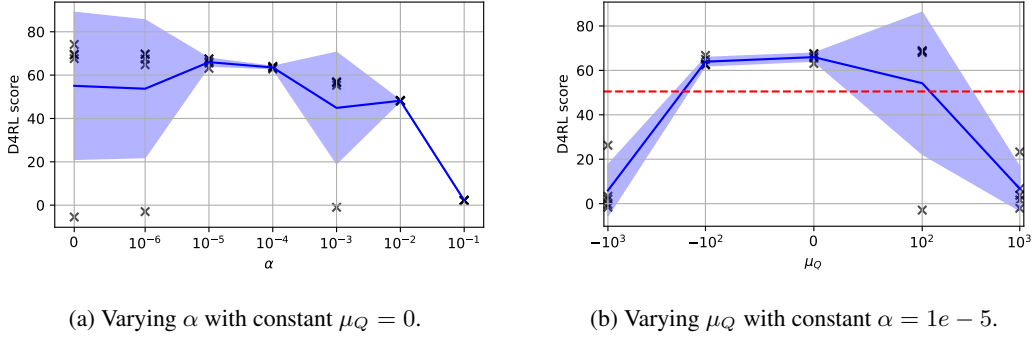


Figure 3: Hyperparameter study on halfcheetah-medium, with 5 seeds per configuration and shaded sample standard deviation. The red line in Fig. 3b corresponds to the average performance after carrying out RP as described in Section 3.

6 Discussion and Limitations

We have presented a Bayesian formulation of the offline RL continuous control problem through conservative value priors. We demonstrated the qualitatively desirable behaviour of our approach on a classic control problem and strong quantitative performance on robotics tasks with larger datasets and complex environments through a simple modification to standard online RL algorithms, outperforming comparable prior model-free methods while remaining highly competitive with more sophisticated and algorithmically complex methods. On small datasets where exact inference is feasible, our method provides consistent epistemic uncertainty quantification. However, since uncertainty quantification is decoupled from decision-making in our framework, we can scale to larger problems without strictly requiring it. By unifying Bayesian decision theory, uncertainty quantification and the challenge of avoiding out-of-distribution actions in offline RL, we establish a foundation for a principled but scalable probabilistic approach to decision-making in offline RL.

Finally, we summarise the limitations of our approach. While in their current form the algorithms provided could also be applied to stochastic environments, our analysis and empirical evaluation is currently restricted to deterministic MDPs. Further, the choice of a suitable conservative prior, along with associated hyperparameters, remains problem-specific and may be challenging in some contexts. The exact inference version of CVP does not scale well to large datasets and, as it relies on hand-crafted kernels, can be ill-suited to high-dimensional environments. While GP-based approximate inference methods [Titsias, 2009] could in principle be applied in our setting to tackle scalability issues, it remains unclear whether they are sufficient for accurately representing a critic function. One limitation towards the applicability of the deep learning, scalable variant of CVP is finding an appropriate regularisation strength α without ground-truth evaluation. Preliminary results suggest that it may be possible to derive reasonable α from training diagnostics (see Appendix E), but further work is necessary to verify this claim’s robustness. We also observe relatively high variance in performance on certain tasks, such as hopper-medium, which is however not uncommon in offline RL algorithms due to the challenging nature of the task. Finally, while building scalable Bayesian uncertainty estimates is in principle possible from the probabilistic formulation we have proposed, deep off-policy evaluation itself is difficult [Fu et al., 2021], so assigning accurate but scalable uncertainty measures is likely to be a challenging task, albeit one that has huge potential for impact to enhance the practical applicability of offline RL.

Acknowledgements

FV was supported by a Department of Computing PhD scholarship and AAF holds a UKRI Turing AI Fellowship (grant no. EP/V025449/1).

References

- Shipra Agrawal and Navin Goyal. Near-optimal regret bounds for thompson sampling. *Journal of the ACM (JACM)*, 64(5):1–24, 2017.
- Gaon An, Seungyong Moon, Jang-Hyun Kim, and Hyun Oh Song. Uncertainty-based offline reinforcement learning with diversified q-ensemble. *Advances in Neural Information Processing Systems*, 34:7436–7447, 2021.
- Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pages 30–37. Morgan Kaufmann, San Francisco (CA), 1995. ISBN 978-1-55860-377-6. doi: <https://doi.org/10.1016/B978-1-55860-377-6.50013-X>.
- David R Burt, Sebastian W Ober, Adrià Garriga-Alonso, and Mark van der Wilk. Understanding variational inference in function-space. *arXiv preprint arXiv:2011.09421*, 2020.
- Ching-An Cheng, Tengyang Xie, Nan Jiang, and Alekh Agarwal. Adversarially trained actor critic for offline reinforcement learning. In *International Conference on Machine Learning*, pages 3852–3878. PMLR, 2022.
- Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. In *Conference on Robot Learning*, pages 885–897. PMLR, 2020.
- Thomas Degris, Martha White, and Richard Sutton. Off-policy actor-critic. *Proceedings of the 29th International Conference on Machine Learning, ICML 2012*, 1, 05 2012.
- Francesco D' Angelo and Vincent Fortuin. Repulsive deep ensembles are Bayesian. In *Advances in Neural Information Processing Systems*, volume 34, pages 3451–3465. Curran Associates, Inc., 2021.
- Michael O’Gordon Duff. *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes*. University of Massachusetts Amherst, 2002.
- Yaakov Engel, Shie Mannor, and Ron Meir. Reinforcement learning with Gaussian processes. In *Proceedings of the 22nd International Conference on Machine Learning*, pages 201–208, 2005.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- Justin Fu, Mohammad Norouzi, Ofir Nachum, George Tucker, Ziyu Wang, Alexander Novikov, Mengjiao Yang, Michael R Zhang, Yutian Chen, Aviral Kumar, et al. Benchmarks for deep off-policy evaluation. *International Conference on Learning Representations*, 2021.
- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. *Advances in Neural Information Processing Systems*, 34:20132–20145, 2021.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International Conference on Machine Learning*, pages 1587–1596. PMLR, 2018.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *International Conference on Machine Learning*, pages 2052–2062. PMLR, 2019.
- Kamyar Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic? Estimating uncertainties for offline rl through ensembles, and why their independence matters. *Advances in Neural Information Processing Systems*, 35:18267–18281, 2022.

- Omer Gottesman, Fredrik Johansson, Matthieu Komorowski, Aldo Faisal, David Sontag, Finale Doshi-Velez, and Leo Anthony Celi. Guidelines for reinforcement learning in healthcare. *Nature Medicine*, 25(1):16–18, 2019.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018a.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018b.
- Danijar Hafner, Dustin Tran, Timothy Lillicrap, Alex Irpan, and James Davidson. Noise contrastive priors for functional uncertainty. In *Uncertainty in Artificial Intelligence*, pages 905–914. PMLR, 2020.
- James Hensman, Nicolò Fusi, and Neil D. Lawrence. Gaussian processes for big data. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, UAI’13, page 282–290, Arlington, Virginia, USA, 2013. AUAI Press.
- Zhiyu Huang, Jingda Wu, and Chen Lv. Efficient deep reinforcement learning with imitative expert priors for autonomous driving. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pages 651–673. PMLR, 2018.
- Alex Kendall, Jeffrey Hawke, David Janz, Przemyslaw Mazur, Daniele Reda, John-Mark Allen, Vinh-Dieu Lam, Alex Bewley, and Amar Shah. Learning to drive in a day. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8248–8254. IEEE, 2019.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Matthieu Komorowski, Leo A Celi, Omar Badawi, Anthony C Gordon, and A Aldo Faisal. The artificial intelligence clinician learns optimal treatment strategies for sepsis in intensive care. *Nature Medicine*, 24(11):1716–1720, 2018.
- Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. Offline reinforcement learning with fisher divergence critic regularization. In *International Conference on Machine Learning*, pages 5774–5783. PMLR, 2021.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. In *International Conference on Learning Representations*, 2022.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in Neural Information Processing Systems*, 30, 2017.
- Jiafei Lyu, Xiaoteng Ma, Xiu Li, and Zongqing Lu. Mildly conservative q-learning for offline reinforcement learning. In *Advances in Neural Information Processing Systems*, 2022.
- Chao Ma and José Miguel Hernández-Lobato. Functional variational inference based on stochastic process generators. *Advances in Neural Information Processing Systems*, 34:21795–21807, 2021.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.

- Andrew William Moore. Efficient memory-based learning for robot control. Technical report, University of Cambridge, 1990.
- Ashvin Nair, Abhishek Gupta, Murtaza Dalal, and Sergey Levine. Awac: Accelerating online reinforcement learning with offline datasets. *arXiv preprint arXiv:2006.09359*, 2020.
- Sebastian W Ober, Carl E Rasmussen, and Mark van der Wilk. The promises and pitfalls of deep kernel learning. In *Uncertainty in Artificial Intelligence*, pages 1206–1216. PMLR, 2021.
- Ian Osband, John Aslanides, and Albin Cassirer. Randomized prior functions for deep reinforcement learning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, David Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *Advances in Neural Information Processing Systems*, 32, 2019.
- Brendan O’Donoghue, Ian Osband, Remi Munos, and Volodymyr Mnih. The uncertainty bellman equation and exploration. In *International conference on machine learning*, pages 3836–3845, 2018.
- Pascal Poupart, Nikos Vlassis, Jesse Hoey, and Kevin Regan. An analytic solution to discrete bayesian reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 697–704, 2006.
- Carl Edward Rasmussen, Christopher KI Williams, et al. *Gaussian processes for machine learning*, volume 1. Springer, 2006.
- Christian P Robert et al. *The Bayesian choice: from decision-theoretic foundations to computational implementation*, volume 2. Springer, 2007.
- Samarth Sinha, Ajay Mandlekar, and Animesh Garg. S4rl: Surprisingly simple self-supervision for offline reinforcement learning in robotics. In *Conference on Robot Learning*, pages 907–917. PMLR, 2022.
- Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- Denis Tarasov, Alexander Nikulin, Dmitry Akimov, Vladislav Kurenkov, and Sergey Kolesnikov. CORL: Research-oriented deep offline reinforcement learning library. In *3rd Offline RL Workshop: Offline RL as a “Launchpad”*, 2022. URL <https://openreview.net/forum?id=SyAS49bBcv>.
- Denis Tarasov, Vladislav Kurenkov, Alexander Nikulin, and Sergey Kolesnikov. Revisiting the minimalist approach to offline reinforcement learning. *Advances in Neural Information Processing Systems*, 36, 2024.
- Michalis Titsias. Variational learning of inducing variables in sparse Gaussian processes. In *Artificial Intelligence and Statistics*, pages 567–574. PMLR, 2009.
- Ahmed Touati, Harsh Satija, Joshua Romoff, Joelle Pineau, and Pascal Vincent. Randomized value functions via multiplicative normalizing flows. In *Uncertainty in Artificial Intelligence*, pages 422–432. PMLR, 2020.
- Filippo Valdettaro and A. Aldo Faisal. Towards offline reinforcement learning with pessimistic value priors. In *Epistemic Uncertainty in Artificial Intelligence*, pages 89–100, Cham, 2024. Springer Nature Switzerland. ISBN 978-3-031-57963-9.
- Joost van Amersfoort, Lewis Smith, Andrew Jesson, Oscar Key, and Yarin Gal. On feature collapse and deep kernel learning for single forward pass uncertainty. *arXiv preprint arXiv:2102.11409*, 2021.
- Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P Xing. Deep kernel learning. In *Artificial Intelligence and Statistics*, pages 370–378. PMLR, 2016.

- Yifan Wu, George Tucker, and Ofir Nachum. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Teng Xiao and Donglin Wang. A general offline reinforcement learning framework for interactive recommendation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 4512–4520, 2021.
- Huan Xu and Shie Mannor. The robustness-performance tradeoff in Markov decision processes. *Advances in Neural Information Processing Systems*, 19, 2006.
- Kai Yang, Jian Tao, Jiafei Lyu, and Xiu Li. Exploration and anti-exploration with distributional random network distillation. In *International Conference on Machine Learning*, pages 56397–56421. PMLR, 2024.
- Wenxuan Zhou, Sujay Bajracharya, and David Held. Plas: Latent action space for offline reinforcement learning. In *Conference on Robot Learning*, pages 1719–1735. PMLR, 2021.

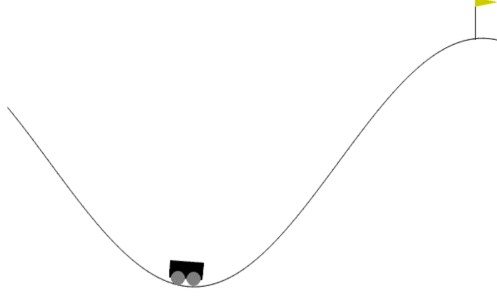


Figure 4: Mountain car environment. The agent’s goal is to push the car starting from the bottom of the valley past the yellow flag on the right.

A Mountain Car Details

The mountain car [Moore, 1990] tackled in Section 4.2 is a deterministic classic control environment. The state-space is two dimensional (position and velocity) and the action space is continuous and one-dimensional $[-1, 1]$. The objective is to strategically apply impulses to a car so that it can escape the valley it starts and build up speed to reach the goal at the top of the mountain, visualised in Fig. 4. If the agent reaches the goal (which is to get the position coordinate to cross past the 0.45 threshold) it receives a reward of 1 and the episode terminates. Otherwise, the reward at each timestep is 0.

The exact-inference CVP agent we train in Section 4.2 has an actor architecture of two hidden layers with 256 neurons each. We train it for 500 gradient steps with Adam optimiser and learning rate of $1e - 3$. In the gradient computation for the optimisation step, in analogy to target networks in standard RL [Mnih et al., 2015, Haarnoja et al., 2018a, Fujimoto et al., 2018], we do not pass gradients through the terms used to evaluate action-values at the next-states $Q^\pi(s', \pi(s'))$, which we treat as stop-gradients for this computation. We use a prior on the action-value function with zero mean, unit variance and RBF kernels, with position, velocity and action lengthscales chosen by hand to be 0.02, 0.008 and 0.5 respectively. Due to the deterministic nature of rewards, we choose a small observation noise variance (the ε in Eq. 3) of $1e-5$.

The expert agent from which the dataset is gathered is taken from the publicly available model at this url: <https://huggingface.co/sb3/sac-MountainCarContinuous-v0/tree/main>.

The TD3 agent is trained for 50,000 training steps using full-batch gradient descent with the same actor network architecture as the CVP agent, an equivalent critic network architecture and default TD3-specific hyperparameters (target critic soft update coefficient 0.005, policy noise 0.2 clipped at 0.5 and delayed policy updates every 2 critic updates) and ADAM optimisers with learning rate of $3e-4$ for both actor and critic. We find, as expected from naive offline application of TD3, that both actor and critic training losses diverge due to the absence of regularisation.

B Neural Network CVP Loss

Here we provide full derivations to show that the loss in Eq. 9 is of the appropriate form for \hat{Q} to regress to the posterior mean for independent state-action kernels as introduced in Section 4.1. We start by considering a finite state-action space, with state space of size $|\mathcal{S}|$ and action space of size $|\mathcal{A}|$, and independent Gaussian priors over Q with mean μ_Q and state-action dependent variance $\sigma_Q^2(s, a)$. With this setup, the following claim holds:

Claim 1. *Using the definition*

$$L_B(Q) = \sum_{(s,a,r,s') \in \mathcal{D}} (Q(s, a) - r - \gamma Q(s', \pi(s')))^2, \quad (12)$$

the loss objective

$$L(Q) = L_B(Q) + \alpha \mathbb{E}_{(s,a) \sim p_{\text{pseudo}}(s,a)} (Q(s, a) - \mu_Q)^2 \quad (13)$$

with $p_{\text{pseudo}}(s, a) = \frac{1}{\alpha} \frac{\sigma_r^2}{\sigma_Q^2(s, a)}$ is minimised when $Q(s, a)$ are the posterior mean.

Proof. Recall that the statistical model from Section 4.1 with observation noise variance σ_r^2 implies a posterior of the form

$$p(Q|\mathcal{D}) \propto p(Q)p(\mathcal{D}|Q) \quad (14)$$

$$\propto \exp \left\{ \sum_{s \in \mathcal{S}, a \in \mathcal{A}} -\frac{(Q(s, a) - \mu_Q)^2}{2\sigma_Q^2(s, a)} \right\} \exp \left\{ -\frac{1}{2\sigma_r^2} \sum_{(s_i, a_i, r_i, s'_i) \in \mathcal{D}} (Q(s_i, a_i) - r_i - \gamma Q(s'_i, \pi(s'_i)))^2 \right\} \quad (15)$$

$$\propto \exp \left\{ -\frac{1}{2} \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \frac{(Q(s, a) - \mu_Q)^2}{\sigma_Q^2(s, a)} - \frac{1}{2\sigma_r^2} \sum_{(s_i, a_i, r_i, s'_i) \in \mathcal{D}} (Q(s_i, a_i) - r_i - \gamma Q(s'_i, \pi(s'_i)))^2 \right\} \quad (16)$$

$$\propto \exp \left\{ -\frac{1}{2\sigma_r^2} \left(\sum_{s \in \mathcal{S}, a \in \mathcal{A}} \sigma_r^2 \frac{(Q(s, a) - \mu_Q)^2}{\sigma_Q^2(s, a)} + \sum_{(s_i, a_i, r_i, s'_i) \in \mathcal{D}} (Q(s_i, a_i) - r_i - \gamma Q(s'_i, \pi(s'_i)))^2 \right) \right\}. \quad (17)$$

Maximising $p(Q|\mathcal{D})$ is equivalent to minimising $-\log p(Q|\mathcal{D})$, so the MAP estimate of Q will minimise

$$\sum_{(s_i, a_i, r_i, s'_i) \in \mathcal{D}} (Q(s_i, a_i) - r_i - \gamma Q(s'_i, \pi(s'_i)))^2 + \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \frac{\sigma_r^2}{\sigma_Q^2(s, a)} (Q(s, a) - \mu_Q)^2. \quad (18)$$

We notice the first term is exactly L_B . Next, we rewrite the sum as an expectation over p_{pseudo} as given in Claim 1 so that the above becomes

$$L_B(Q) + \sum_{s \in \mathcal{S}, a \in \mathcal{A}} p_{\text{pseudo}}(s, a) \underbrace{\frac{1}{p_{\text{pseudo}}(s, a)} \frac{\sigma_r^2}{\sigma_Q^2(s, a)}}_{\alpha} (Q(s, a) - \mu_Q)^2. \quad (19)$$

$$= L_B(Q) + \alpha \mathbb{E}_{(s, a) \sim p_{\text{pseudo}}(s, a)} (Q(s, a) - \mu_Q)^2 \quad (20)$$

after substituting α as defined in the claim. We have thus shown that the MAP Q -values will be recovered by minimising Eq. 13, which is quadratic in Q . The last step to prove Claim 1 is to notice that for Gaussian variables MAP and posterior mean coincide, since the mode of a Gaussian is equal to its mean. \square

Having established that in the tabular case with finite state-action space the loss in Eq. 13 provably results in posterior mean Q -values, we consider increasing the state-action space to an intractably large, approximately continuous, number of state. Standard arguments [Mnih et al., 2015] suggest replacing the tabular $Q(s, a)$ with a parametrised $\hat{Q}_\theta(s, a)$, with training objective given by the same loss but replacing $Q(s, a)$ with $\hat{Q}_\theta(s, a)$, leading directly to the loss in Eq. 9.

Notice that the relationship

$$\sigma_Q^2(s, a) = \frac{\sigma_r^2}{\alpha} \frac{1}{p_{\text{pseudo}}(s, a)} \quad (21)$$

implies that the regularisation strength constant α can be interpreted as being inversely proportional to the prior variance, as claimed in Section 5.1.

Finally, we comment on the specific choice of σ_Q^2 and, correspondingly, p_{pseudo} . A natural baseline, as done when applying Alg. 1 in practice, is to take a constant $\sigma_Q^2(s, a) = \sigma_Q^2$. However, this would entail using a uniform $p_{\text{pseudo}}(s, a) = \frac{1}{|\mathcal{S}||\mathcal{A}|}$, which is problematic for large or high-dimensional state-action spaces. As a consequence, as described in Section 5.1, we depart from using such a uniform p_{pseudo} in our experiments. Despite this modification, a Bayesian interpretation remains valid, but with a state-dependent prior variance given by Eq. 21. In the context of our decision-making task, this corresponds to setting narrower priors, thus stronger regularisation, on those state-actions that are most relevant to decision-making.

C Continuous maze

Here we compare the results of applying the exact-inference based Alg.1 and the deep-learning critic variant as described in Section 5.1 and Alg. 2, with two base online algorithms, one learning a deterministic (Twin Delayed DDPG [Fujimoto et al., 2018]), and one a stochastic (Soft Actor-Critic [Haarnoja et al., 2018a,b]) policy, with offline CVP-based algorithms called CVP-TD3 and CVP-SAC respectively. We choose to compare these approaches on a structured toy environment, with easily interpretable qualitative behaviour.

Continuous maze We consider an MDP in a continuous-space maze-like environment where only a part of the state-action space is adequately covered by the dataset. The state-space is $\mathcal{S} = \mathbb{R}^2$ and the action space is $\mathcal{A} = [-1, 1]^2$, visualised in Fig. 5a. The agent deterministically moves in state space by adding the action vector to its current state vector. It receives a reward of 1 upon reaching the goal region $\{(x, y) \in \mathbb{R}^2 : 2 < x < 3, 0 < y < 1\}$ and 0 otherwise, with episodes terminating at the goal.

The static dataset we consider covers a limited part of the environment’s state-space, and is formed by individual (s, a, r, s') transitions rather than full episodic traces. The datasets contains transitions starting from 100 uniformly spaced states in the region $\{(x, y) \in \mathbb{R}^2 : 2 < x < 3, 0 < y < 1\} \cup \{(x, y) \in \mathbb{R}^2 : 0 < x < 3, 1 < y < 2\}$. All actions in the dataset are in either the x or y directions and have size 0.6, with the action observed at each state being the one that makes progress towards the goal while remaining in the in-support region (see Fig. 5a).

Exact inference We first apply Alg. 1 to this toy dataset, with RBF kernels with length-scale 0.25 for both state and actions and an actor parametrised by an MLP with two hidden layers of 256 neurons trained for 1000 steps. In Fig 5b, we observe that the resulting policy reaches the goal while remaining in the in-distribution data regions. The values learned are displayed in Fig. 5c, where we observe that only the in-distribution regions are assigned high values, which gradually decrease with number of steps required to reach the goal (as expected due to the discount factor $\gamma = 0.9$). Similarly, the value uncertainty displayed in Fig. 5d is high in the OOD regions and low where the agent can confidently reach the goal while remaining in-distribution. The thin lines of slightly decreased expected value in some in-distribution regions of Fig. 5c are due to the continuous nature of the actor. In these regions that bridge a sharp policy change, continuity implies that the actor must employ actions that are slightly different to those on either side of the transition region, thus resulting in a lower covariance with the neighbouring states and lower posterior value.

CVP-SAC We start by considering CVP with SAC as base algorithm. For consistency with the GP implementation, we use prior mean $\mu_Q = 0$, and employ $\alpha = 0.001$ for both methods. We implement our algorithm by modifying base implementations provided by Tarasov et al. [2022].

Visualisations of the learned policy and values from applying Alg. 2 are shown in Fig. 5e. The scalable version produces a policy that avoids the no-data region while reaching the goal, with a value estimate consistent with the exact-inference case in the in-distribution region. Since the pseudo-data distribution regularises OOD actions rather than states, OOD states don’t receive a signal to align with the prior. In contrast, Fig. 5f shows that applying the base SAC algorithm naively to the offline data (equivalent to setting $\alpha = 0$) leads to poor policies that do not remain in the support of the data and wander into regions of state-space that the agent has no knowledge of, leading to unsatisfactory offline policies. We examine the learned values in action space for this toy experiment further in Appendix C.

TD3-SAC. Next, we report analogous results for the the TD3-based CVP-TD3. In Fig. 5g, we observe how the CVP-TD3 policy successfully solves the offline RL task and learns similar values to those learned with exact inference CVP in the in-distribution region, whereas in Fig. 5h we see that the unregularised naive application of TD3 fails in this instance. This closely mirrors the behaviour shown in Section C with SAC.

Visualising CVP in action space. In Fig.6, we observe the effect of the conservative prior directly in action-space, where we visualise the action-values after training in the *action* space at state $(2.5, 1.1)$ just above the goal region. There is an action $(0, -0.6)$ in the dataset at this state that leads directly to the goal, which we expect the agent to follow. For all three methods, the neighbourhood

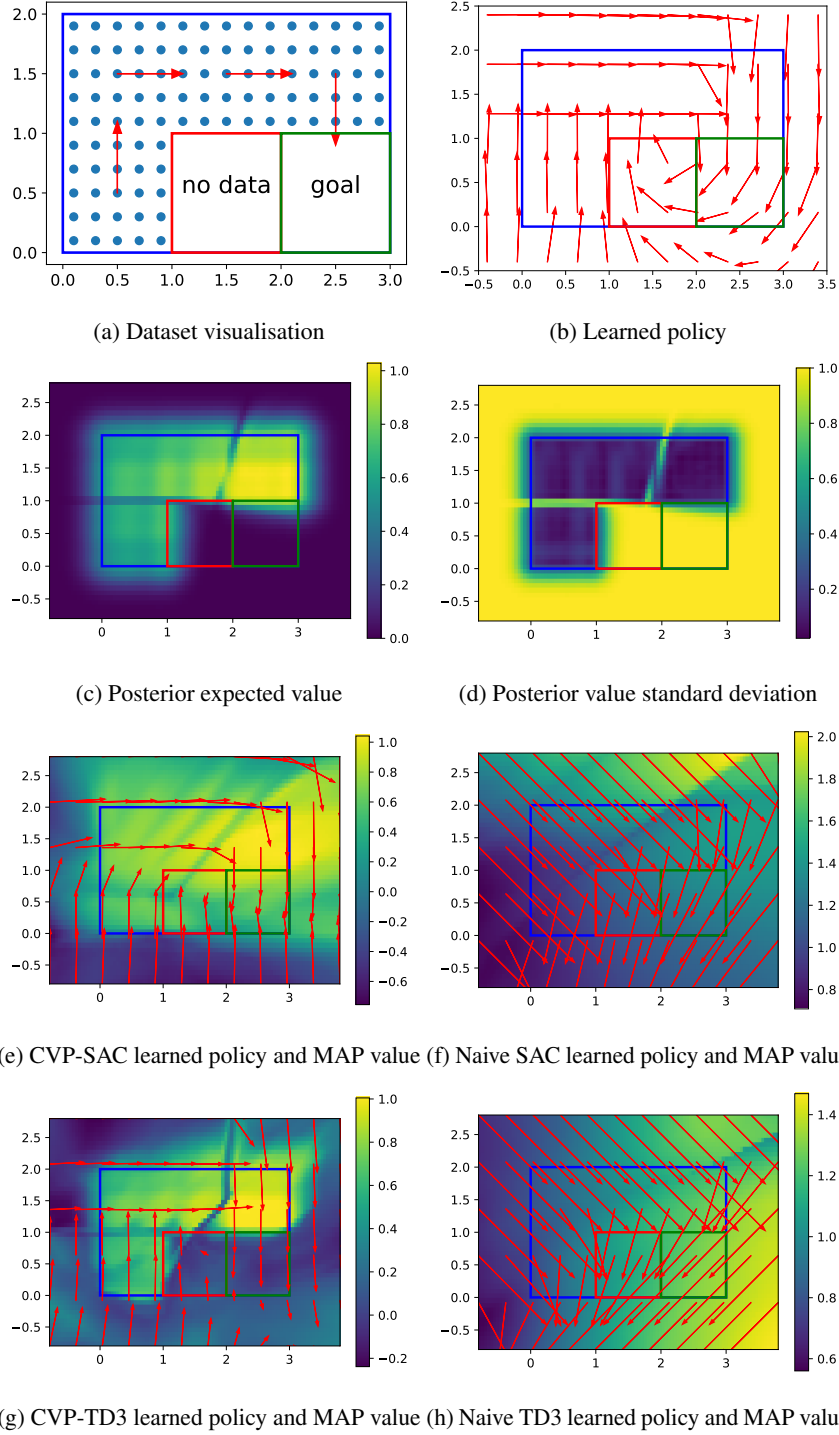


Figure 5: Toy environment for navigation task. The agent receives a reward of 1 for transitioning into the goal region (where the episode terminates) and 0 otherwise. The dataset consists of steps of size 0.6 in the cardinal direction that leads towards the goal while remaining in the supported (blue) region. Learned policies and value functions are visualised as arrows and colour-maps respectively. The length of the arrows corresponds to the size of the step taken.

around the observed action that leads to the goal is regularised as desired, and the algorithms correctly learn to choose the action similar to the one in the dataset that leads to the goal. As can be seen in

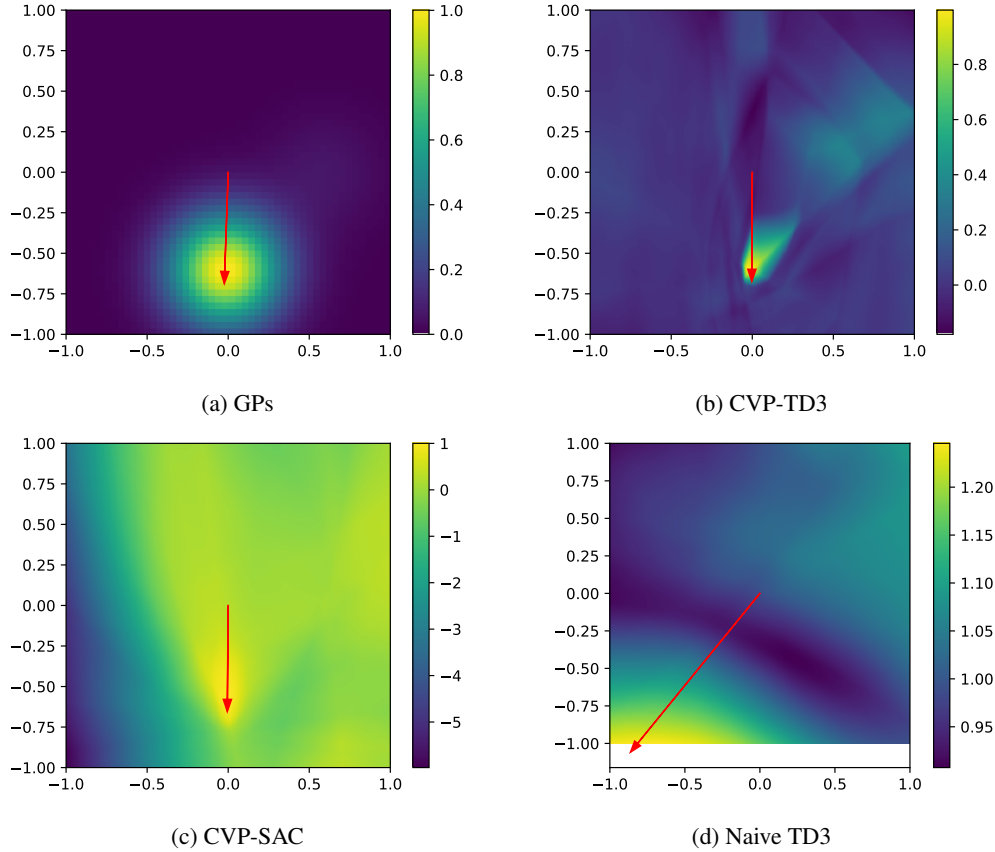


Figure 6: Learned action-values at state $(2.5, 1.1)$, located just above the goal. The red arrow shows the learned action. The dataset contains an action $(0, -0.6)$ leading to the goal at this state.

Fig. 6a, due to the exact nature of the inference employed, the GPs smoothly regularise the whole action space around the observed action. Similarly, the actions that maximise value in Figs. 6b and 6c correctly correspond to following the observed action leading to the goal. While the extrapolation in regions far from the action the agent is considering can vary significantly across methods, this is because our approximate methods must focus on regularising those actions most towards the actor’s decision making, and we do observe that values learned close to the actor’s action are similar across methods, therefore ensuring that the appropriate action is learned. In contrast, Fig. 6d shows how a naive application of TD3 (setting $\alpha = 0$ in TD3-CVP) causes the actor to choose an unsupported action, where the increase in value is not supported by observed actions but rather entirely caused by unwarranted extrapolation.

D D4RL Implementation and Compute Details

We base our implementation of CVP-SAC, which we train on the D4RL benchmarks in Section 5.2 on the SAC backbone of the CQL implementation in the CORL library [Tarasov et al., 2022], and the results presented do not require any additional hyperparameter tuning to the base algorithm beyond changing the optimiser of the entropy coefficient loss to Adam. The hyperparameter relevant to the base SAC algorithm, common to all tasks, are summarised in Table 2.

Next, we address the CVP-specific hyperparameters. We report the prior mean used for each task in Table 3 (either choosing μ_Q or carrying out reward preprocessing as explained in Section 3) and the choice of α for each task. Other methods that carry out critic value regularisation, such as CQL [Kumar et al., 2020] and MCQ [Lyu et al., 2022], both sample 10 actions at which to regularise, so we take the parameter $n = 10$ in Alg. 2 and do not tune it further.

Hyperparameter	Value
Batch size	256
Discount factor	0.99
Training steps	1e6
Actor hidden layers	3
Actor hidden size	256
Actor optimiser	Adam
Actor learning rate	3e-5
Critic hidden layers	3
Critic hidden size	256
Critic optimiser	Adam
Critic learning rate	3e-4
Target network update rate	0.005
Entropy coefficient optimiser	SGD
Entropy coefficient learning rate	3e-5

Table 2: Base SAC hyperparameters for the D4RL experiments.

Table 3: CVP hyperparameter settings for the D4RL experiments. RP in the μ_Q column refers to using a mean of 0 after carrying out reward preprocessing as described in Section 3.

Task	α	μ_Q
halfcheetah-random	1e-4	0
hopper-random	1e-4	RP
walker2d-random	1e-2	0
halfcheetah-medium	1e-5	0
hopper-medium	5e-4	0
walker2d-medium	5e-3	RP
halfcheetah-medium-replay	5e-5	0
hopper-medium-replay	5e-4	RP
walker2d-medium-replay	5e-3	RP
halfcheetah-medium-expert	1e-2	0
hopper-medium-expert	5e-3	-10
walker2d-medium-expert	1e-2	0
halfcheetah-expert	1e-2	0
hopper-expert	1e-1	0
walker2d-expert	1e-2	0

Hyperparameter tuning. We started by fixing $\mu_Q = 0$ as default and tuning α per-task. The heuristic described in Appendix E guided the α tuning, giving an indication of whether the next α to try should be greater, if the training was unstable due to regularisation being too small as can be inferred from training diagnostics, or smaller, to check if the regularisation can be decreased further, with possible performance benefits. We then tried RP with these values of α . For some tasks, such as hopper-medium-expert, we observed stable training diagnostics but higher variance in performance, suggesting that further tuning of μ_Q or prior mean selection could be beneficial, and used a different value for μ_Q if it had a beneficial effect on performance.

Compute resources. The experiments were carried out on standard general-purpose workstations provided by the authors’ affiliated institution, rather than on dedicated high-performance computing systems. As the networks employed are small (three-layer MLPs with 256 neurons) the experiments were ran on a mixture of GPU and CPU-only machines, depending on availability. A typical runtime for a full training run would be of approximately 10.5 hours on a machine with a NVIDIA GeForce GTX 1050 Ti GPU and a Intel(R) Core(TM) i7-10700 CPU @ 2.90GHz CPU. This includes extensive evaluation in simulation and logging of training diagnostics during training, and we did not explicitly optimise our code for performance.

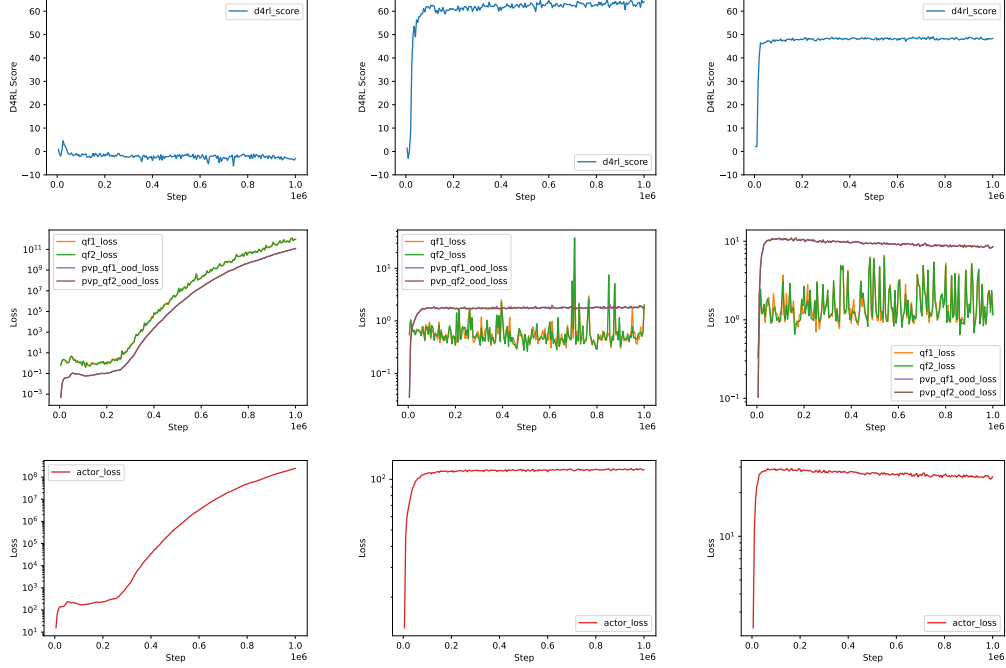


Figure 7: Comparison of training diagnostics and performance for different α values. The left column displays results for $\alpha = 1e - 6$, the middle column $\alpha = 1e - 4$ and the right column for $\alpha = 1e - 2$. qf_loss is the standard RL Bellman loss of the 2 critics used whereas ood_loss is the CVP loss term.

E Choosing regularisation strength from training diagnostics

Here we show how we can empirically guide the choice of regularisation strength, α in Alg. 2, off of training diagnostics - without relying on evaluations on the ground-truth environment. This points to the possibility of automated tuning of this hyperparameter, which we leave for future work.

We take as reference the halfcheetah-medium task. As mentioned in Section 5.3, we notice that reducing regularisation strength can be beneficial up to a certain critical point where training becomes unstable. A weak regularisation leads to inconsistent results, but the maximum achieved returns can be very high, even higher than that achieved by the α that gives the best average across 5 seeds. Meanwhile, regularisation that is too strong slowly degrades performance.

Further investigating the training diagnostics reveals that the failure modes causing this suboptimal performance arise from very different mechanisms. We find that strong regularisation generally leads to stable training (low losses) but poor fitting to the value function whereas small regularisation has greater potential for better performance but at the risk of very high instability. We display relevant losses for sample runs with different α that show this effect in Fig. 7. This strongly suggests that the optimal regularisation strength should be as small as possible while still ensuring stable training and losses. Indeed, there is strong evidence that the optimal value for α might depend on initialisation, as some seeds with the lowest regularisation setting still achieve strong performance with stable convergence, albeit more rarely than when a slightly higher α is employed, and some seeds with higher regularisation can still occasionally exhibit unstable training. This leads to the natural hypothesis that an automated tuning of α would be able to further increase the performance of CVP-SAC beyond what can be achieved by employing a constant α during training, which we leave for future work. We found the pattern of choosing the lowest α before training instability arises as generally leading to good ground-truth evaluation performance in the tasks considered, but further experimentation is required to determine the robustness of this empirical observation.

We show training diagnostics for exemplar runs in Fig. 7. We plot a diverging exemplar run for when α is too small ($\alpha = 10^{-6}$), where training is unstable and we observe training losses and Q-values

diverge in a similar way as they would without regularisation. When α is too large ($\alpha = 10^{-2}$), training is stable but the returns achieved are suboptimal. This is in line with the notion that excessive regularisation can harm performance, and the ideal α should be the small but large enough to prevent the training from diverging. We remark that, in this example, this can be identified directly by the training diagnostics, without requiring evaluation on the ground truth environment.

F Additional D4RL results

We compare here our results to two other model-free offline RL algorithms, that however require significant algorithmic extensions to standard off-policy algorithms achieve a performance similar to that of CVP-SAC, in contrast to those presented in the main Table 1. In this section, we also summarise the nature of these additional components. Overall, CVP-SAC achieves 97% and 102% of the scores attained by MCQ and ATAC, respectively, on the D4RL locomotion tasks.

MCQ comparison. We compare to the scores obtained by the Mildly Conservative Q Learning (MCQ) algorithm [Lyu et al., 2022] in Table 4. Similarly to our method, MCQ includes a L2 regularisation term in the critic loss, but the value against which the critic is regularised is learned during training, requiring the introduction of a generative model for the behaviour policy, thus significantly increasing the complexity of the algorithm. We see that CVP is able to recover 97% of MCQ’s average performance without requiring such a generative model.

Table 4: Comparison between MCQ [Lyu et al., 2022] and CVP results on D4RL benchmarks, indicating the percentage of MCQ’s score that achieved by CVP-SAC, without having to explicitly model the behaviour policy. The MCQ results are taken from Lyu et al. [2022].

Task	MCQ	CVP	Percentage achieved
halfcheetah-random	28.5	32.5	114.0%
hopper-random	31.8	27.5	86.5%
walker2d-random	17.0	6.9	40.6%
halfcheetah-medium	64.3	66.0	102.6%
hopper-medium	78.4	72.3	92.2%
walker2d-medium	91.0	84.1	92.4%
halfcheetah-medium-replay	56.8	57.2	100.7%
hopper-medium-replay	101.6	100.5	98.9%
walker2d-medium-replay	91.3	81.5	89.3%
halfcheetah-medium-expert	87.5	95.4	109.0%
hopper-medium-expert	111.2	103.8	93.3%
walker2d-medium-expert	114.2	109.3	95.7%
halfcheetah-expert	96.2	95.1	98.9%
hopper-expert	111.4	110.7	99.4%
walker2d-expert	107.2	110.0	102.6%
Average	79.2	76.9	97.0%

ATAC comparison. We report an analogous set of results comparing to the Adversarially Trained Actor Critic (ATAC) for offline RL algorithm [Cheng et al., 2022] in Table 5. Compared to the other baseline methods in Section 5.2, ATAC includes additional ingredients to ensure the adversarial approach they use has stable training. For example, they introduce updates on different timescales for actor and critic as well as a term in the critic loss, novel to offline RL, that combines a weighted double Q loss and a residual algorithm loss [Baird, 1995], which they found to be crucial for performance. Neither of these optimisation-related ingredients is strictly related to offline RL, and we did not implement these into CVP in our experiments. Nevertheless, we find that CVP still is able to achieve slightly better performance (102%) than ATAC in these tasks.

Table 5: Comparison between ATAC [Cheng et al., 2022] and CVP results on D4RL benchmarks, indicating the percentage of ATAC’s score achieved by CVP-SAC. The ATAC results are taken from Cheng et al. [2022], where the expert datasets are not reported.

Task	ATAC	CVP	Percentage of score achieved
halfcheetah-random	3.9	32.5	833.3%
hopper-random	17.5	27.5	157.1%
walker2d-random	6.8	6.9	101.5%
halfcheetah-medium	53.3	66.0	123.8%
hopper-medium	85.6	72.3	84.4%
walker2d-medium	89.6	84.1	93.9%
halfcheetah-medium-replay	48	57.2	119.2%
hopper-medium-replay	102.5	100.5	98.0%
walker2d-medium-replay	92.5	81.5	85.2%
halfcheetah-medium-expert	94.8	95.4	98.8%
hopper-medium-expert	111.9	103.8	92.7%
walker2d-medium-expert	114.2	109.3	95.7%
Average	68.4	69.7	102.0%