OFFLINE REINFORCEMENT LEARNING THROUGH TRAJECTORY CLUSTERING AND LOWER BOUND PENALISATION

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we propose a new framework for value regularisation in offline reinforcement learning (RL). While most previous methods evade explicit out-of-distribution (OOD) region identification due to its difficulty, our method explicitly identifies the OOD region, which can be non-convex depending on datasets, via a newly proposed trajectory clustering-based behaviour cloning algorithm. With the obtained explicit OOD region, we then define a Bellman-type operator pushing the value in the OOD region to a tight lower bound while operating normally in the in-distribution region. The value function with this operator can be used for policy acquisition in various ways. Empirical results on multiple offline RL benchmarks show that our method yields the state-of-the-art performance.

1 Introduction

Offline reinforcement learning (RL) has gained significant attention from the RL community due to its efficiency in cost and safety. Unlike conventional RL, where an agent learns an optimal policy through interactions with the environment, in offline RL, environmental interactions are not done. Instead, a set \mathcal{D} of trajectories is provided and the agent searches for a competent policy using only the samples in \mathcal{D} . Despite its attractiveness, there exists a critical hurdle to offline RL: overestimation of Q values of the critic network in the out-of-distribution (OOD) action region (Fujimoto et al., 2019). Since the agent cannot actually perform an overestimated action and gain correction signals from the environment, the extrapolation error will be corrected directly.

Critic penalisation is one of the main offline approaches to handle this value overestimation problem in the OOD region (Kumar et al., 2020; Lyu et al., 2022; Mao et al., 2023). By penalising the critic values of penalized OOD actions, this approach nudges the agent to select in-distribution (ID) actions with higher critic values rather than OOD ones. Various offline algorithms have been developed for critic penalisation with various penalisation terms. Due to the difficulty of identifying the OOD region itself, these methods rely on indirect measures to construct a penalisation term, e.g., difference-based penalization (Kumar et al., 2020), importance sampling-based integration (Mao et al., 2023). However, these indirect penalisation methods have their own shortcomings for the cost of evading direct OOD region identification. For example, SVR (Mao et al., 2023) uses simple Gaussian behaviour modelling to compute the required importance sampling ratio. When the dataset has two equally-strong modes with a reasonable distance, the learned Gaussian-modelled behaviour density will place a large density value at the centre of the two modes on which the actual density is very low, and the algorithm will yield significantly degraded performance.

To overcome such limitations resulting from not identifying the precise OOD region, in this paper, we proposes a novel value regularisation algorithm that explicitly identifies the multi-modal OOD region and penalises the critic values of the OOD region with a newly-derived lower bound tighter than previous approaches. Empirical evaluations of our method on the D4RL benchmark (Fu et al., 2020) show that our method yields high-performing policies from various offline RL datasets. The main contributions of our work are:

• We introduce a criterion for identifying whether an action is OOD based on its likelihood.

- Connecting trajectory clustering in offline RL to task identification in meta RL, we propose an algorithm that can learn mixture Gaussian approximations to the behaviour policy.
- We develop a new value regulariser that regresses the critic values toward a tight theoretical lower bound of the optimal action-value function.

2 BACKGROUND

Notation For the list of notations used in this paper and their meanings, refer to Appendix Sec. A.

Markov Decision Process An RL problem is formulated as a Markov Decision Process (MDP), which is defined as a 6-tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma, \rho_0 \rangle$, where $\mathcal{S} \subseteq \mathbb{R}^{d_s}$ is the state space, $\mathcal{A} \subseteq \mathbb{R}^{d_a}$ is the action space, $P \colon \mathcal{S} \times \mathcal{A} \to \mathcal{P}(\mathcal{S})$ is the transition dynamics, $R \colon \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathcal{P}(\mathbb{R})$ is the reward function, $\gamma \in [0,1]$ is the discount factor, and $\rho_0 \in \mathcal{P}(\mathcal{S})$ is the initial state distribution. We will assume that the support of R(s,a,s') is bounded above by r_{\max} and bounded below by r_{\min} for all $s,s' \in \mathcal{S}$ and $a \in \mathcal{A}$.

Value Functions Given a policy π , the Bellman operator \mathcal{T}^{π} on $L^{\infty}(\mathcal{S} \times \mathcal{A})$ is defined by the following equation:

$$(\mathcal{T}^{\pi}Q)(s,a) = \mathbb{E}_{s' \sim P(s,a)} \left[\mathbb{E}_{r \sim R(s,a,s')}[r] \right] + \mathbb{E}_{s' \sim P(s,a)} \left[\mathbb{E}_{a' \sim \pi(s')} \left[Q(s',a') \right] \right].$$

Then, the action-value function (or Q-function) $Q^{\pi} \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is defined as the unique fixed point of \mathcal{T}^{π} , and the state-value function $V^{\pi} \colon \mathcal{S} \to \mathbb{R}$ is given by $V^{\pi}(s) = \mathbb{E}_{a \sim \pi(s)} \left[Q^{\pi}(s, a) \right]$. The objective of RL is to find an optimal policy π^* such that $V^{\pi^*} \succeq V^{\pi}$ for any policy π .

Offline Reinforcement Learning For offline RL, interactions with the environment is prohibited, and the agent has to learn a policy from a given dataset \mathcal{D} of trajectories. Throughout this paper, we will assume that each trajectory $\tau \in \mathcal{D}$ is sampled with a uni-modal behaviour policy $\beta \in \{\beta_0, \beta_1, \beta_2, \dots, \beta_{K-1}\}$, where the candidate set $\mathcal{B} = \{\beta_0, \beta_1, \beta_2, \dots, \beta_{K-1}\}$ is fixed but unknown to the agent.

3 MOTIVATION

Critic penalization or value regularisation penalises the Q-values for OOD actions, while minimizing the temporal difference error for in-distribution (ID) actions. We may formulate it with the following equation

$$\min_{Q} \mathbb{E}_{(s,a) \sim \mathcal{D}} \left[\left(Q(s,a) - \mathcal{T}^{\pi} Q(s,a) \right)^{2} \right] + \mathfrak{R},$$

where \Re is a regularizer. A crucial requirement of the regularizer is that it should be able to discriminate between ID and OOD actions since we only want to penalise the values of OOD actions. One of the first approaches was to set the regulariser as (Kumar et al., 2020)

$$\mathfrak{R} = \mathbb{E}_{s \sim \mathcal{D}, a \sim \mu}[Q(s, a)] - \mathbb{E}_{s \sim \mathcal{D}, a \sim \beta}[Q(s, a)],$$

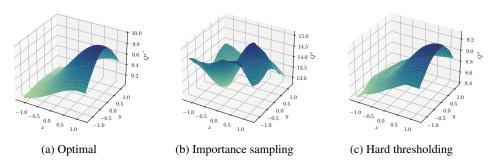


Figure 1: The Q-values on the plane spanned by e_1 and e_2 estimated by each method, i.e., $Q(xe_1 + ye_2)$. Due to the high variance of importance sampling ratios, the importance sampling method fails to approximate the optimal Q-function accurately.

where β is the behaviour policy and μ is some distribution that satisfies the condition $\operatorname{supp} \mu \subseteq \operatorname{supp} \beta$ (Kumar et al., 2020). While minimising the Q values for OOD actions sampled by μ , they simultaneously maximised the Q values for ID actions sampled from β to compensate for overpenalisation. However, as Mao et al. (2023) points out, this approach has two shortcomings: (i) the requirement $\operatorname{supp} \mu \subseteq \operatorname{supp} \beta$ may not hold in general; and (ii) if the dataset contains a large portion of suboptimal actions, their Q values would be overestimated. To address these issues, they proposed an importance sampling (IS)-based method that utilises the following regulariser:

$$\mathfrak{R}_{\mathrm{IS}} = \mathbb{E}_{s \sim \mathcal{D}, a \sim \mu} [(Q(s, a) - Q_{\mathrm{targ}}(s, a))^2] - \mathbb{E}_{s \sim \mathcal{D}, a \sim \beta} \left[\frac{\mu(a \mid s)}{\beta(a \mid s)} (Q(s, a) - Q_{\mathrm{targ}}(s, a))^2 \right],$$

where μ is a probability distribution supported on the entire action space and $Q_{\rm targ}$ is a regulariser target, which they set to $r_{\rm min}/(1-\gamma)$ for all $s\in\mathcal{S}, a\in\mathcal{A}$. Since the two terms cancel each other on supp β , $\mathfrak{R}_{\rm IS}$ is equivalent to $\mathbb{E}_{s\sim\mathcal{D},a\sim\mu}\left[\mathbf{1}_{\mathcal{A}\setminus\text{supp}\,\beta}(Q(s,a)-Q_{\rm targ}(s,a))^2\right]$, which corresponds to the goal of penalising the Q values of OOD actions.

A significant drawback of $\mathfrak{R}_{\mathrm{IS}}$ is that IS ratios are known to have high variance, especially for high-dimensional spaces. Consider a simple single-state infinite-horizon MDP with a six-dimensional action space, and an offline RL dataset of size $1\,000\,000\,$ sampled from a behaviour policy $\mathcal{N}(\mathbf{0}, \mathbf{I}_6)$. Suppose the optimal action is $\mathbf{a}^* = \mathbf{e}_1$. Then the IS ratio between $\mu = \mathcal{N}(\mathbf{a}^*, 0.04\mathbf{I}_6)$ and β of the samples in the dataset ranges from 1.88×10^{-225} to 1.93×10^4 . As demonstrated in Figure 1b and Table 1, the IS method yields an inaccurate Q value estimation and a suboptimal policy due to this severe fluctuation of IS ratios.

Table 1: The discounted return of the policies learned with each method. IS stands for importance sampling and HT stands for hard thresholding.

Optimal	IS	НТ	
10	2.44 ± 0.83	9.72 ± 0.14	

To overcome these limitations of the previous value regularization methods, we here propose to explicitly identify the set of OOD actions OOD(s) for each state $s \in \mathcal{S}$, and set the regulariser to zero for ID actions. Such hard thresholding (HT) allows a more stable training process, resulting in a more accurate Q value estimations and better-performing policies, as seen in Fig. 1c and Table 1.

4 PROPOSED METHOD

This section is structured as follows. In Section 4.1, we first discuss how we can compute the set OOD(s). We then propose a new lower bound of Q^{π^*} and show its effectiveness as a penalisation target in Section 4.2. Finally, we provide a practical offline RL algorithm in Section 4.3.

4.1 IDENTIFYING THE OUT-OF-DISTRIBUTION ACTION SET

Likelihood is the most natural way to measure how OOD a particular sample is. However, choosing the threshold value is not trivial. For blunt distributions, we should use a lower threshold value, whereas for sharp distributions, we can choose a higher value. We propose a systematic method of setting the threshold value by adopting the concept of *highest density region* (HDR; Hyndman 1996), which is basically a generalisation of a confidence interval to multivariate random variables.¹

Definition 1 (Hyndman 1996). Let f(X) be the pdf of a random variable X. Then, the $100(1-\alpha)\%$ highest density region (HDR) is the subset $\mathcal{R}(f_{\alpha})$ of the sample space of X such that $\mathcal{R}(f_{\alpha}) = \{x : f(x) \ge f_{\alpha}\}$, where $f_{\alpha} = \sup\{y : \mathbb{P}(X \in \mathcal{R}(y)) \ge 1 - \alpha\}$.

In the following subsections, we discuss how to compute the HDR under different assumptions.

4.1.1 Homogeneous Datasets

We first discuss the case when the offline dataset \mathcal{D} is homogeneous, that is, it was generated from a single uni-modal behaviour policy β . Then, we may obtain a fairly accurate Gaussian approximation

 $^{^1}$ We provide a diagram (Figure 5) showing the $100(1-\alpha)$ % HDR of a normal distribution on page 27 to aid the understanding of the concept of a HDR.

 $\hat{\beta}$ of β through behaviour cloning. Let $\mu \colon \mathcal{S} \to \mathbb{R}^{d_a}$ and $\Sigma \colon \mathcal{S} \to \mathbb{R}^{d_a \times d_a}$ be the mean and covariance matrix functions of $\hat{\beta}$, respectively. Assuming $\Sigma(s)$ is positive definite for all $s \in \mathcal{S}$, the $100(1-\alpha)\%$ HDR has the following closed-form representation (Proposition 4 in Appendix):

$$\mathcal{R}_{\hat{\beta}}(f_{\alpha}; s) = \left\{ \mathbf{x} \in \mathbb{R}^{d_{a}} : A_{\hat{\beta}}(\mathbf{x}; s) \leq F_{\chi_{d}^{2}}^{-1}(1 - \alpha) \right\},\,$$

where $F_{\chi^2_{d_a}}^{-1}$ is the inverse cumulative distribution function of a chi-squared random variable with d_a degrees of freedom and $A_{\hat{\beta}}(\mathbf{x};s) = (\mathbf{x} - \boldsymbol{\mu}(s))^{\top} \boldsymbol{\Sigma}(s)^{-1} (\mathbf{x} - \boldsymbol{\mu}(s))$. Choosing an appropriate value of $0 < \alpha < 1$, we can define $\mathrm{OOD}(s)$ as

$$OOD(s) = A \setminus \mathcal{R}_{\hat{\beta}}(f_{\alpha}; s). \tag{1}$$

4.1.2 HETEROGENEOUS DATASETS

The definition of OOD(s) for a homogeneous dataset given in (1) can be generalised to the heterogeneous case as $OOD(s) = \mathcal{A} \setminus \left(\bigcup_{\beta \in \mathcal{B}} \mathcal{R}_{\beta}(f_{\alpha}; s)\right)$ for \mathcal{B} , where \mathcal{B} is the behaviour policy candidate set. If we could identify and isolate all of the trajectories in the dataset sampled from a particular behaviour policy $\beta \in \mathcal{B}$, then obtaining an estimation $\hat{\beta}$ of β is straightforward by applying a behaviour cloning algorithm on those isolated trajectories. Then, with the estimated $\hat{\mathcal{B}} = \{\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_{K-1}\}$, we could compute OOD(s) for each state s as above. Therefore, in the rest of this section, we will propose how to cluster the trajectories. Note that the proposed clustering algorithm is useful not only for value regularization here but also for other offline real-world data analysis.

Our key idea is that the trajectory clustering problem closely resembles the task inference problem in meta RL. For each policy π , there is a corresponding Markov reward process (MRP) $\mathcal{M}^{\pi} = \langle \mathcal{S}, P^{\pi}, R^{\pi}, \gamma \rangle$, where for all $s, s' \in \mathcal{S}$ and $r \in \mathbb{R}$, the transition probability function $P^{\pi} \colon \mathcal{S} \to \mathcal{P}(\mathcal{S})$ and the reward probability function $R^{\pi} \colon \mathcal{S} \times \mathcal{S} \to \mathcal{P}(\mathbb{R})$ are defined by the equations

$$P^{\pi}(s'\mid s) = \mathbb{E}_{a\sim\pi(\cdot\mid s)}\left[P(s'\mid s, a)\right],\tag{2}$$

$$R^{\pi}(r \mid s, s') = \mathbb{E}_{a \sim \pi(\cdot \mid s)} \left[R(r \mid s, a, s') \right], \tag{3}$$

respectively. Since the dataset \mathcal{D} can then be viewed as a collection of trajectories, where each trajectory is sampled from one of the MRPs \mathcal{M}^{β_0} , \mathcal{M}^{β_1} , \mathcal{M}^{β_2} , ..., $\mathcal{M}^{\beta_{K-1}}$, trajectory clustering task can be viewed as an MRP inference problem. As this formulation is almost equivalent to the MDP inference problem setting in meta RL, we infer the MRP instead of the MDP and apply a technique similar to variational Bayes-adaptive deep RL (variBAD; Zintgraf et al. 2021).

Our goal is to infer the behaviour policy index given a trajectory. To achieve this objective, we represent the index as a discrete latent variable m supported on $[K] = \{0, 1, \dots, K-1\}$ and write

$$P^{\beta_m}(s) \approx P(s; m),$$
 $R^{\beta_m}(s, s') \approx R(s, s'; m),$ $\beta_m(s) \approx \beta(s; m),$

for all $s,s' \in \mathcal{S}$, sharing P,R, and β across trajectories. The marginal pdf of a trajectory $\tau_{:T} = (s_0,a_0,r_0,s_1,a_1,r_1,s_2,a_2,r_2,\ldots,s_{T-1},a_{T-1},r_{T-1},s_T)$ is

$$p(\tau_{:T}) = \rho_0(s_0) \sum_{m=0}^{K-1} p(m) \prod_{t=0}^{T-1} P(s_{t+1} \mid s_t; m) \beta(a_t \mid s_t; m) R(r_t \mid s_t, s_{t+1}; m), \tag{4}$$

where p(m) is the prior distribution on m. Modelling P, R, and β with neural networks parametrised by θ results in a loss that depends on θ . However, the multi-modality of (4) causes gradient-based optimisation algorithms to produce sub-optimal solutions. We circumvent this issue by introducing amortised inference network q_{ϕ} that takes a variable-length action-less trajectory $\tilde{\tau}_{:t} = (s_0, r_0, s_1, r_1, s_2, r_2, \dots, s_{t-1}, r_{t-1}, s_t)$ as an input and outputs a distribution in $\mathcal{P}_d([K])$. Instead of maximising (4), we maximise the evidence lower bound (ELBO), which can be written by the following equation (Proposition 5):

$$ELBO_{\theta,\phi}(\tau;t) = -D_{KL}(q_{\phi}(\tilde{\tau}_{:t}) \parallel p) + \sum_{i=0}^{T-1} \mathbb{E}_{m \sim q_{\phi}(\tilde{\tau}_{:t})} \left[\log R_{\theta}(r_i \mid s_i, s_{i+1}; m) \right]$$
 (5)

$$+ \sum_{i=0}^{T-1} \mathbb{E}_{m \sim q_{\phi}(\tilde{\tau}_{:t})} \left[\log \beta_{\theta}(a_i \mid s_i; m) \right] + \sum_{i=0}^{T-1} \mathbb{E}_{m \sim q_{\phi}(\tilde{\tau}_{:t})} \left[\log P_{\theta}(s_{i+1} \mid s_i; m) \right].$$

Encoder $\begin{array}{c}
h_{t-1} \\
s_t \longrightarrow q_{\phi} \\
r_t \longrightarrow q_{\phi}
\end{array}
\longrightarrow b_t = q_{\phi}(\cdot \mid \tilde{\tau}_{:t}) \sim m \xrightarrow{s_{i+1}} S_{i} \longrightarrow \beta_{\theta} \longrightarrow a_i$

Figure 2: Overview of our architecture. We also provided a diagram of the variBAD architecture in Figure 6 for comparison.

The first term $\log \rho_0(s_0)$ in (9) can be omitted because it is constant with respect to θ and ϕ . The final objective for trajectory clustering is to maximise

$$\mathbb{E}_{\tau \sim \mathcal{D}} \left[\frac{1}{T_{\tau}} \sum_{t=0}^{T_{\tau}-1} \text{ELBO}_{\theta,\phi}(\tau;t) \right], \tag{6}$$

where T_{τ} is the length of the trajectory τ sampled from the dataset. An overview of our clustering algorithm is given in Figure 2.

After we finish training, we compute the behaviour policy estimations and cluster assignments according to the equations $\hat{\beta}_i = \beta_{\theta}(\cdot\,;i)$ and $\mathbb{A}(s) = \arg\max_m q_{\phi}(m\mid \tilde{\tau}(s))$, respectively, for each $i\in[K]$ and $s\in\mathcal{D}$, where $\tilde{\tau}(s)$ is the action-less trajectory containing the state s.

4.2 LOWER BOUND PENALISATION

In this section, we derive a new lower bound on the value function. As we did in the previous section, we start with the case where the offline dataset $\mathcal D$ is generated from a single behaviour policy β . The ideal penalisation method would be to use $Q^{\pi^*}(s,a)$ as a target, where π^* is the optimal policy, but the value of Q^{π^*} is inaccessible. So we aim to use a lower bound instead. In order to compute a lower bound, we first need to make some assumptions on the regularity of P and V^{β} .

Assumption 1. There is $K_P > 0$ such that for all $s \in \mathcal{S}$ and $a, a' \in \mathcal{A}$, $W_1(P(s, a), P(s, a')) < K_P || a - a' ||$, where $W_1(P, Q)$ is the Wasserstein distance of order 1 between two probability distributions $P, Q \in \mathcal{P}(\mathcal{S})$.

Assumption 2. The value function of the behaviour policy β is K_V -Lipschitz, that is, for all $s, s' \in \mathcal{S}$, $|V^{\beta}(s) - V^{\beta}(s')| < K_V ||s - s'||$.

Then, we can obtain a lower bound of Q^{π^*} with these assumptions.

Proposition 1. Define $Q^{\mathrm{LB}}_{\beta} \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ by the equation

$$Q_{\beta}^{\text{LB}}(s, a) = \max \left\{ V^{\beta}(s) - r_{\text{max}} + r_{\text{min}} - \gamma K_V K_P \mathbb{E}_{a' \sim \beta(s)} \left[\|a - a'\| \right], \frac{r_{\text{min}}}{1 - \gamma} \right\}. \tag{7}$$

For any policy $\pi \colon \mathcal{S} \to \mathcal{P}(\mathcal{A})$ such that $V^{\pi} \succeq V^{\beta}$, $Q^{\pi} \succeq Q^{\mathrm{LB}}_{\beta}$.

Proof. See page 14.
$$\Box$$

Note that this lower bound is tighter than the previous bound $r_{\min}/(1-\gamma)$. The lower bound allows us to define the penalised Bellman optimality operator $\mathcal{T}^{\pi}_{\beta}$ for policy π by the equation

$$(\mathcal{T}_{\beta}^*Q)(s,a) = \begin{cases} Q_{\beta}^{\mathrm{LB}}(s,a) & \text{if } a \in \mathrm{OOD}(s), \\ (\mathcal{T}^*Q)(s,a) & \text{otherwise,} \end{cases}$$

where \mathcal{T}^* is the Bellman optimality operator defined as

$$(\mathcal{T}^*Q)(s,a) = \mathbb{E}_{s' \sim P(s,a), r \sim R(s,a,s')} \left[r + \gamma \sup_{a' \in \mathcal{A}} Q(s',a') \right].$$

We can show that through repeated application of \mathcal{T}^*_{β} , it is possible to obtain a deterministic policy $\pi^*_{\beta} \colon \mathcal{S} \to \mathcal{A}$ that is optimal among the policies whose action for each state $s \in \mathcal{S}$ does not lie in OOD(s).

Theorem 2. Any initial bounded real-valued function on $S \times A$ can converge to a unique fixed point Q_{β}^* by repeatedly applying \mathcal{T}_{β}^* . Suppose for each $s \in S$,

$$Q^{\beta}(s, a_s) \ge \mathbb{E}_{a \sim \beta(s)} \left[Q^{\beta}(s, a) \right]$$

for some $a_s \in \mathcal{A} \setminus \mathrm{OOD}(s)$. If there exists a deterministic policy $\pi_{\beta}^* \colon \mathcal{S} \to \mathcal{A}$ that is optimal under the constraint $\pi(s) \notin \mathrm{OOD}(s)$ for all $s \in \mathcal{S}$, then $\pi_{\beta}^*(s) = \arg\max_{a \in \mathcal{A}} Q_{\beta}^*(s, a)$ for all $s \in \mathcal{S}$.

Proof. See page 19.
$$\Box$$

Now, the penalised Bellman optimality operator can easily be generalised to the heterogeneous dataset case with the set \mathcal{B} of behaviour policy candidates and the set $\mathcal{V}(s)$ of valid behaviour policies given a state $s \in \mathcal{S}$.

$$(\mathcal{T}_{\mathcal{B}}^*Q)(s,a) = \begin{cases} Q_{\mathcal{B}}^{\mathrm{LB}}(s,a) & \text{if } a \in \mathrm{OOD}(s), \\ (\mathcal{T}^*Q)(s,a) & \text{otherwise,} \end{cases}$$

where $Q_{\mathcal{B}}^{\mathrm{LB}} \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is defined as $Q_{\mathcal{B}}^{\mathrm{LB}}(s,a) = \max_{\beta \in \mathcal{V}(s)} Q_{\beta}^{\mathrm{LB}}(s,a)$ for each $s \in \mathcal{S}$ and $a \in \mathcal{A}$. We can prove a similar performance guarantee for the policy obtained by repeatedly applying $\mathcal{T}_{\mathcal{B}}^*$.

Theorem 3. Any initial bounded real-valued function on $S \times A$ can converge to a unique fixed point Q_B^* by repeatedly applying \mathcal{T}_B^* . Suppose for each $\beta \in \mathcal{B}$ and $s \in \mathcal{S}$,

$$Q^{\beta}(s, a_s^{\beta}) \ge \mathbb{E}_{a \sim \beta(s)} \left[Q^{\beta}(s, a) \right]$$

for some $a_s^\beta \in \mathcal{A} \setminus \mathrm{OOD}(s)$. If there exists a deterministic policy $\pi_\mathcal{B}^* \colon \mathcal{S} \to \mathcal{A}$ that is optimal under the constraint $\pi(s) \not\in \mathrm{OOD}(s)$ for all $s \in \mathcal{S}$, then $\pi_\mathcal{B}^*(s) = \arg\max_{a \in \mathcal{A}} Q_\mathcal{B}^*(s,a)$ for all $s \in \mathcal{S}$.

Proof. See page 20.
$$\Box$$

4.3 PRACTICAL ALGORITHM

The overall flow of our algorithm is as follows:

- I. Behaviour policy learning. Run the trajectory clustering algorithm to obtain $\hat{\mathcal{B}}$ and $\hat{\mathcal{V}}(s)$. Or if it is known a priori that the dataset is homogeneous, then run a behaviour cloning algorithm to obtain $\hat{\beta}$.
- II. Behaviour value learning. Learn a value function $\hat{V}^{\hat{\beta}}$ for each $\hat{\beta} \in \hat{\mathcal{B}}$ through temporal difference learning.
- III. **Policy learning.** Obtain and apply $\mathcal{T}_{\mathcal{B}}^*$ repeatedly on a randomly initialised Q-function until convergence. Find a policy that maximises the learned Q-function.

This section mainly focuses on the trajectory clustering algorithm of Stage I. Additional details of our algorithm can be found in Section C of Appendix.

The network architecture used for trajectory clustering consists of three parts: the encoder, the latent sampler, and the decoder. The architecture is generally similar to that of variBAD except for a few adaptations. In this section, we will first go over how and why we modified each part. Then, we will propose a simple technique to adaptively set the number of clusters.

The encoder needs to take an action-less trajectory $\tilde{\tau}$ as an input and output the amortised posterior. Since the length of the trajectory may vary from one to another, the network should be capable

of taking variable-length sequence as its input. For that purpose, variBAD utilises gated recurrent units (GRU; Cho et al. 2014). GRUs and other recurrent neural network variants suffer from the vanishing gradient problem (Bengio et al., 1994), which hampers their ability to process long sequences. Truncated backpropagation through time (Williams & Peng, 1990) can mitigate the phenomena to a certain extent, but we instead adopt the state space model architecture that is recently gaining interest in the area of sequence modelling (Gu et al., 2020; 2021; 2022; Gu & Dao, 2023; Dao & Gu, 2024). In particular, we use the S5 layer (Smith et al., 2023), which is simple and computationally efficient.

The second modification was made on the way latents are sampled and ELBOs are computed. As the latent variable in the variBAD architecture is continuous, it is impossible to analytically compute the expectation, and hence, the reparametrisation trick (Kingma & Welling, 2014) must be used. Although the latent variable is discrete in our case, exact computation is still inefficient because it requires multiple forward and backward passes through the decoder. We instead utilise the vector quantised-variational autoencoder (VQ-VAE; van den Oord et al. 2017) to approximate the ELBO. Under the VQ-VAE formulation, the amortised posterior q_{ϕ} is modelled as

$$q_{\phi}(m = k \mid \tilde{\tau}_{:t}) = \begin{cases} 1 & \text{if } k = k_{\phi}(\tilde{\tau}_{:t}) \\ 0 & \text{otherwise,} \end{cases}$$

where e_0, e_1, \dots, e_{K-1} are latent embedding vectors and

$$k_{\phi}(\tilde{\tau}_{:t}) = \underset{j \in [K]}{\operatorname{arg\,min}} \|q_{\phi}(\tilde{\tau}_{:t}) - e_k\|_2.$$

Note that for simplicity, we have abused the notation q_{ϕ} to denote both the posterior and the encoder. The gradient flows into the encoder q_{ϕ} via the loss function

$$\ell_{\text{VQ}}(\phi; \tilde{\tau}_{:t}) = \left\| q_{\phi}(\tilde{\tau}_{:t}) - e_{k_{\phi}(\tilde{\tau}_{:t})} \right\|_{2}^{2}$$

and the latent embedding vectors are updated with exponential moving averages.

For the decoder, we use Gaussian distributions with diagonal covariance matrix to represent P_{θ} , β_{θ} , and R_{θ} . Most RL environments have a bounded action space, whereas a Gaussian distribution has unbounded support. To estimate the behaviour policy more accurately, we first normalize the actions between -1 and 1 and apply the inverse hyperbolic tangent function on each dimension of the actions to map them onto \mathbb{R}^{d_a} . Note that we use the mapped actions when learning the critic, that is, the critic function takes $\tanh^{-1}(a)$ instead of a as input. Finally, instead of taking the summation over the entire trajectory in (5) and (6), we adopt the implementation trick of variBAD and randomly subsample N_d transition steps in (5) and N_e ELBO terms in (6). To conclude, the loss function for the trajectory clustering algorithm is

$$\ell_{\text{TC}}(\theta, \phi; \tau) = \frac{1}{N_e N_d} \sum_{t \in \mathcal{I}_e} \sum_{i \in \mathcal{I}_d} A_{\theta}(s_i, a_i, r_i, s_{i+1}; e_{k_{\phi}(\tilde{\tau}:t)}) + \lambda_{\text{VQ}} \ell_{\text{VQ}}(\phi; \tilde{\tau}_{:t}),$$

where

$$A_{\theta}(s, a, r, s'; m) = \log \beta_{\theta}(a \mid s; m) + \lambda_T \log P_{\theta}(s' \mid s; m) + \lambda_R \log R_{\theta}(r \mid s, s'; m), \tag{8}$$

 \mathcal{I}_e and \mathcal{I}_d are sets of indices sampled uniformly at random with replacement from $[T_{\tau}]$ with sizes N_e and N_d , respectively, and $\lambda_{\text{VQ}}, \lambda_T, \lambda_R$ are tunable hyperparameters.

Choosing the right number of clusters is crucial for high performance in most clustering algorithms. To alleviate the burden of hyperparameter tuning, we adopt a two-phase training paradigm. During the first phase of the paradigm, we set the codebook size to be sufficiently large. After completing the first phase, we compute the cluster assignments for each state in the dataset. If the number of states assigned to a particular cluster does not exceed a certain threshold, we remove the corresponding code from the VQ-VAE codebook. The training is resumed with the remaining codebook. This way, we could adaptively determine the number of clusters without needing to perform an exhaustive hyperparameter search.

Table 2: Average normalised scores on the D4RL benchmark. Note that "ha" means halfcheetah, "ho" means hopper, "wa" means walker2d, "m" means medium, "r" means replay, "ra" means random, and "e" means expert.

Dataset	BC	TD3BC	BCQ	BEAR	CQL	IQL	MCQ	SVR	Ours
ha-ra	2.6	11.0	2.2	2.3	17.5	13.1	28.5	27.2	27.0 ± 1.1
ho-ra	4.1	8.5	7.8	3.9	7.9	7.9	31.8	31.0	31.5 ± 0.2
wa-ra	1.2	1.6	4.9	12.8	5.1	5.4	17.0	2.2	16.6 ± 7.9
ha-m	42.0	48.3	46.6	43.0	47.0	47.4	64.3	60.5	63.5 ± 0.9
ho-m	56.2	59.3	59.4	51.8	53.0	66.2	78.4	103.5	102.8 ± 0.4
wa-m	71.0	83.7	71.8	-0.2	73.3	78.3	91.0	92.4	94.1 ± 2.2
ha-m-r	36.4	44.6	42.2	36.3	45.5	44.2	56.8	52.5	52.2 ± 0.8
ho-m-r	21.8	60.9	60.9	52.2	88.7	94.7	101.6	103.7	102.2 ± 1.1
wa-m-r	24.9	81.8	57.0	7.0	81.8	73.8	91.3	95.6	95.4 ± 19.2
ha-m-e	59.6	90.7	95.4	46.0	75.6	86.7	87.5	94.2	
ho-m-e	51.7	98.0	106.9	50.6	105.6	91.5	111.2	111.2	112.4 ± 1.1
wa-m-e	101.2	110.1	107.7	22.1	107.9	109.6	114.2	109.3	108.3 ± 0.7
ha-e	88.2	81.7	92.7	92.9	96.3	95.0	96.2	96.1	96.6 ± 0.9
ho-e	110.9	107.8	109.0	54.6	96.5	109.4	111.4	111.1	112.7 ± 0.9
wa-e	107.7	110.2	106.3	106.6	108.5	109.9	107.2	110.0	113.4 ± 0.5
Average	52.3	67.5	64.5	38.8	67.3	68.9	79.2	80.0	

Table 3: The performance of SVR and our method on the custom heterogeneous dataset.

Algorithm	Length	Return		
SVR	8.00 ± 0.00	4.05 ± 0.01		
Ours	436.3 ± 32.1	4062.0 ± 24.5		

5 EXPERIMENTS

5.1 RESULTS ON THE D4RL BENCHMARK

In order to evaluate how well our algorithm perform on various offline RL tasks, we tested our method on the D4RL (Fu et al., 2020) benchmark. We compared it with existing offline RL methods such as BC (Pomerleau, 1988), TD3+BC (Fujimoto & Gu, 2021), CQL (Kumar et al., 2020), IQL (Kostrikov et al., 2022), MCQ (Lyu et al., 2022), and SVR (Mao et al., 2023). We trained our method with five different seeds to obtain five different policies and sampled ten trajectories with each of them. We report the average and standard deviation of the fifty normalized scores in Table 2. The results show that our algorithm can successfully learn high-performing policies from most datasets, while attaining state-of-the-art scores on some of them.

5.2 Experiments on a Heterogeneous Dataset

Although D4RL datasets such as "hopper-medium-expert-v2" were sampled with more than one behaviour policies, the action distributions are actually unimodal on most states due to the state distribution being so different between the two behaviour policies. Figure 3 presents a visualisation of the entire and initial state distributions of the "hopper-medium-expert-v2" dataset where we have used the uniform manifold approximation and projection (UMAP; McInnes & Healy 2018) technique for dimension reduction. We can see that expert and medium states are clearly separated, except for the initial states.

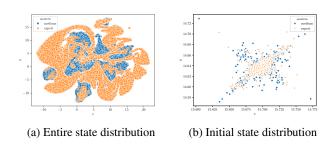


Figure 3: The UMAP of the states in the "hopper-medium-expert-v2" dataset.

To demonstrate the effectiveness of our trajectory clustering algorithm, we created a custom dataset with drastically different initial state behaviours using the "Hopper-v5" environment provided by Gymnasium library. Half of the samples in the dataset were sampled from an expert policy, and the other half was sampled from a policy that tripped over within eight timesteps. Table 3 demonstrates that our method can effectively classify the two datasets and learn an optimal policy from a heterogeneous dataset.

5.3 ANALYSIS ON THE TRAJECTORY CLUSTERING ALGORITHM

In meta reinforcement learning settings, each MDP has independent transition and reward dynamics, so they must be modelled in order to infer the MDP from trajectories. Under our formulation, on the other hand, transition and reward dynamics of each MRP are correlated with each other through the policy as we can see from (2) and (3). Although this implies that we may identify the MRP solely through modelling the behaviour policy, we hypothesized that modelling transition and reward dynamics can provide meaningful auxiliary information leading to better clustering performance. Therefore, we compared the performance of our algorithm under four different configurations $(\lambda_T, \lambda_R) \in \{(1,1), (1,0), (0,1), (0,0)\},\$ where λ_T and λ_R are the weights for transition

Table 4: The impact of hyperparameters λ_T and λ_R on the average performance of our trajectory clustering algorithm evaluated on six custom D4RL datasets. The performance is measured in terms of adjusted rand index (ARI) and normalized mutual information score (NMI).

λ_R	λ_T	ARI	NMI
0	0	0.98 ± 0.07	0.98 ± 0.06
0	1	0.91 ± 0.21	0.92 ± 0.17
1	0	0.99 ± 0.02	0.98 ± 0.02
1	1	0.86 ± 0.27	0.87 ± 0.24

and reward models defined in (8). To evaluate the accuracy of our trajectory clustering algorithm, we created custom D4RL datasets by concatenating random, medium, and expert datasets. The mean and standard deviation of adjusted rand indices (ARI; Hubert & Arabie 1985) and normalised mutual information scores (NMI) for each configuration over 5 different seeds are reported in Table 4. We can see that the configuration $(\lambda_T, \lambda_R) = (1,0)$ performs the best on average. Unlike s_{i+1} , which is in the vicinity of s_i regardless of the a_i , r_i can vary drastically between policies, making it difficult to model rewards from different policies with a single neural network. We speculate this to be the reason why training a reward model negatively affects the performance of our trajectory clustering algorithm. For experiments on other datasets, refer to Section E.3.

6 CONCLUSION

In this paper, we propose a new value regularisation algorithm for offline RL penalizing their critic values, based on the OOD action set that we were able to explicitly identify. We determine how OOD an action is based on its likelihood, where the threshold is set adaptively according to the shape of the behaviour policy. To enable likelihood analysis for heterogeneous datasets where simple behaviour cloning fails, we introduce a novel trajectory clustering technique based on a meta-learning formulation of the clustering problem. Our method of penalising the critic values for OOD actions by regressing them towards a lower bound of the optimal Q-value function is proven to be effective both theoretically and empirically.

REFERENCES

- Yoshua Bengio, Patrice Y. Simard, and Paolo Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Networks*, 5(2):157–166, 1994.
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL http://github.com/jax-ml/jax.
- Kyunghyun Cho, Bart van Merrienboer, Çaglar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pp. 1724–1734. ACL, 2014.
- Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net, 2024.
- Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: datasets for deep data-driven reinforcement learning. *CoRR*, abs/2004.07219, 2020.
- Scott Fujimoto and Shixiang Shane Gu. A minimalist approach to offline reinforcement learning. In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 20132–20145, 2021.
- Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without exploration. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pp. 2052–2062. PMLR, 2019.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *CoRR*, abs/2312.00752, 2023.
- Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. In *Advances in Neural Information Processing Systems 33:* Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.
- Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 572–585, 2021.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *The Tenth International Conference on Learning Representations, ICLR* 2022, *Virtual Event, April* 25-29, 2022. OpenReview.net, 2022.
- Jonathan Heek, Anselm Levskaya, Avital Oliver, Marvin Ritter, Bertrand Rondepierre, Andreas Steiner, and Marc van Zee. Flax: A neural network library and ecosystem for JAX, 2024. URL http://github.com/google/flax.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.*
- Lawrence Hubert and Phipps Arabie. Comparing partitions. *Journal of Classification*, 2(1):193–218, Dec 1985. ISSN 1432-1343. doi: 10.1007/BF01908075. URL https://doi.org/10.1007/BF01908075.
 - Rob J. Hyndman. Computing and graphing highest density regions. *The American Statistician*, 50 (2):120–126, 1996.

- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In 2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings, 2014.
- Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022.* OpenReview.net, 2022.
- Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline reinforcement learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020.*
- Jiachen Li, Edwin Zhang, Ming Yin, Qinxun Bai, Yu-Xiang Wang, and William Yang Wang. Offline reinforcement learning with closed-form policy improvement operators. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 20485–20528. PMLR, 2023.
- Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings, 2016.
- Jianlan Luo, Perry Dong, Jeffrey Wu, Aviral Kumar, Xinyang Geng, and Sergey Levine. Action-quantized offline reinforcement learning for robotic skill learning. In Conference on Robot Learning, CoRL 2023, 6-9 November 2023, Atlanta, GA, USA, volume 229 of Proceedings of Machine Learning Research, pp. 1348–1361. PMLR, 2023.
- Jiafei Lyu, Xiaoteng Ma, Xiu Li, and Zongqing Lu. Mildly conservative q-learning for offline reinforcement learning. In Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022.
- Yihuan Mao, Chengjie Wu, Xi Chen, Hao Hu, Ji Jiang, Tianze Zhou, Tangjie Lv, Changjie Fan, Zhipeng Hu, Yi Wu, Yujing Hu, and Chongjie Zhang. Stylized offline reinforcement learning: Extracting diverse high-quality behaviors from heterogeneous datasets. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.
- Yixiu Mao, Hongchang Zhang, Chen Chen, Yi Xu, and Xiangyang Ji. Supported value regularization for offline reinforcement learning. In *Advances in Neural Information Processing Systems* 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 16, 2023, 2023.
- Leland McInnes and John Healy. UMAP: uniform manifold approximation and projection for dimension reduction. *CoRR*, abs/1802.03426, 2018.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- Dean Pomerleau. ALVINN: an autonomous land vehicle in a neural network. In *Advances in Neural Information Processing Systems 1, [NIPS Conference, Denver, Colorado, USA, 1988]*, pp. 305–313. Morgan Kaufmann, 1988.
- Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. Simplified state space layers for sequence modeling. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023.* OpenReview.net, 2023.
- Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 2256–2265. JMLR.org, 2015.

Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 6306–6315, 2017.

- Cédric Villani. Cyclical monotonicity and Kantorovich duality, pp. 51–92. Springer Berlin Heidelberg, Berlin, Heidelberg, 2009. ISBN 978-3-540-71050-9.
- Qiang Wang, Yixin Deng, Francisco Roldan Sanchez, Keru Wang, Kevin McGuinness, Noel E. O'Connor, and Stephen J. Redmond. Dataset clustering for improved offline policy learning. *CoRR*, abs/2402.09550, 2024.
- Zhendong Wang, Jonathan J. Hunt, and Mingyuan Zhou. Diffusion policies as an expressive policy class for offline reinforcement learning. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- Ronald J. Williams and Jing Peng. An efficient gradient-based algorithm for on-line training of recurrent network trajectories. *Neural Comput.*, 2(4):490–501, 1990.
- Luisa M. Zintgraf, Sebastian Schulze, Cong Lu, Leo Feng, Maximilian Igl, Kyriacos Shiarlis, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: Variational bayes-adaptive deep RL via meta-learning. J. Mach. Learn. Res., 22:289:1–289:39, 2021.

A NOTATIONS

- 0: a zero vector with dimensionality implied by context
- D_{KL}(P₁ || P₂): the Kullback–Leibler (KL) divergence from a probability distribution P₁ to another probability distribution P₂
- e_i: the *i*-th standard basis of a Euclidean space
- $f(y \mid x)$: the value of the pdf (or pmf) of the distribution f(x) at y, where Y is a set and $f: X \to \mathcal{P}(Y)$
- $f \succeq g$: $f(x) \ge g(x)$ for all $x \in X$, where f and g are real-valued functions defined on a set X
- $f \equiv g$: f(x) = g(x) for all $x \in X$, where f and g are real-valued functions defined on a set X
- I_d : an identity matrix with d rows and d columns
- L^{\infty}(X): the space of bounded real value functions on a set X endowed with the supremum norm
- [N]: the set $\{0, 1, \dots, N-1\}$, where N is an integer
- $\mathcal{N}(\mu, \Sigma)$: a multi-variate Gaussian distribution with mean vector μ and covariance matrix Σ
- $\mathbb{P}(E)$: probability of an event E
- $\mathcal{P}(X)$: family of absolutely continuous probability distributions with finite first moments supported on a subset of X, where $X \subseteq \mathbb{R}^d$
- $\mathcal{P}_d(X)$: the family of discrete distributions supported on a subset of X, where $X \subseteq \mathbb{R}^d$
- supp μ : the support of a probability distribution μ
- $\mathcal{U}(X)$: the uniform distribution on a Borel set $X \subseteq \mathbb{R}^d$ of positive, finite Lebesgue measure
- $W_1(P_1, P_2)$: the Wasserstein distance of order 1 between two probability distributions $P_1, P_2 \in \mathcal{P}(X)$

We also define a clipping function

$$\operatorname{clip}(x; y, z) = \max\{y, \min\{x, z\}\}.$$

The notation can be generalised to dimension-wise clipping, that is, for $\mathbf{x} \in \mathbb{R}^d$ and $y, z \in \mathbb{R}$, the *i*-th coordinate of $\operatorname{clip}(\mathbf{x}; y, z)$ is $\operatorname{clip}(x_i; y, z)$.

B PROOFS

 Proposition 4. Let X be a multivariate Gaussian random variable with mean vector $\mu \in \mathbb{R}^d$ and positive definite covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$. The $100(1 - \alpha)$ % HDR is

$$\mathcal{R}(f_{\alpha}) = \left\{ \mathbf{x} \in \mathbb{R}^{d} : (\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq F_{\chi_{d}^{d}}^{-1} (1 - \alpha) \right\},\,$$

where $F_{\chi^2_d}$ is the cumulative distribution function of a chi-squared random variable with d degrees of freedom.

Proof. Let $\mathbf{Z} = (Z_1, Z_2, \dots, Z_d) = \sqrt{\mathbf{\Sigma}^{-1}} (\mathbf{X} - \boldsymbol{\mu})$. By the change of variables formula,

$$p_{\mathbf{Z}}(\mathbf{z}) = \left| \det(\sqrt{\Sigma}) \right| p_{\mathbf{X}} \left(\boldsymbol{\mu} + \sqrt{\Sigma} \mathbf{z} \right)$$
$$= \det(\Sigma)^{1/2} (2\pi)^{-d/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2} \mathbf{z}^{\top} \mathbf{z} \right)$$
$$= (2\pi)^{-d/2} \exp\left(-\frac{1}{2} \mathbf{z}^{\top} \mathbf{z} \right),$$

where $p_{\mathbf{X}}$ and $p_{\mathbf{Z}}$ are the pdfs of random vectors \mathbf{X} and \mathbf{Z} , respectively. We can see that \mathbf{Z} is a standard normal random vector. Since

$$\mathcal{R}(y) = \left\{ \mathbf{x} \in \mathbb{R}^d : (2\pi)^{-d/2} \det(\mathbf{\Sigma})^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \ge y \right\}$$
$$= \left\{ \mathbf{x} \in \mathbb{R}^d : (\mathbf{x} - \boldsymbol{\mu})^{\top} \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \le -2 \log y + d \log(2\pi) + \log \det(\mathbf{\Sigma}) \right\},$$

we have

$$\mathbb{P}(\mathbf{X} \in \mathcal{R}(y)) = \mathbb{P}\left(\mathbf{Z}^{\top}\mathbf{Z} \le -2\log y + d\log(2\pi) + \log\det(\Sigma)\right)$$
$$= \mathbb{P}\left(\sum_{i=1}^{d} Z_i^2 \le -2\log y + d\log(2\pi) + \log\det(\Sigma)\right).$$

 Z_1, Z_2, \ldots, Z_d are independent, so $\sum_{i=1}^d Z_i^2$ is a chi-squared random variable. This implies $\mathbb{P}(\mathbf{X} \in \mathcal{R}(y)) = F_{\chi_{\mathcal{A}}^2}(-2\log y + d\log(2\pi) + \log\det(\mathbf{\Sigma})).$

$$\mathbb{P}(\mathbf{X} \in \mathcal{R}(y)) \ge 1 - \alpha$$
 if and only if

$$-2\log y + d\log(2\pi) + \log \det(\mathbf{\Sigma}) \ge F_{\chi_{\mathcal{J}}^2}^{-1}(1-\alpha).$$

Therefore.

$$f_{\alpha} = (2\pi)^{d/2} \det(\mathbf{\Sigma})^{1/2} \exp\left(-\frac{1}{2}F_{\chi_d^2}(1-\alpha)\right),$$

which means

$$\mathcal{R}(f_{\alpha}) = \left\{ \mathbf{x} \in \mathbb{R}^d : (\mathbf{x} - \boldsymbol{\mu})^{\top} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \leq F_{\chi_d^2}^{-1} (1 - \alpha) \right\}.$$

Proposition 5. Let m be a discrete latent variable supported on [K] and

$$\tau_{:T} = (s_0, a_0, r_0, s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T)$$

be a trajectory rolled-out according to the following sampling process: $s_0 \sim \rho_0$, $m \sim p$, and for each $t \in [T]$, $s_{t+1} \sim P(s_t; m)$, $a_t \sim \beta(s_t; m)$, and $r_t \sim R(s_t, s_{t+1}; m)$. The marginal pdf can be written as

$$p(\tau_{:T}) = \rho_0(s_0) \sum_{m=0}^{K-1} p(m) \prod_{t=0}^{T-1} P(s_{t+1} \mid s_t ; m) \beta(a_t \mid s_t ; m) R(r_t \mid s_t, s_{t+1} ; m)$$

and for any distribution q on [K],

$$\log p(\tau_{:T}) \ge \log \rho_0(s_0) - D_{\text{KL}}(q \parallel p)$$

$$+ \sum_{t=0}^{T-1} \mathbb{E}_{m \sim q} \left[\log P(s_{t+1} \mid s_t; m) + \log \beta(a_t \mid s_t; m) + \log R(r_t \mid s_t, s_{t+1}; m) \right].$$
(9)

Proof. Let us denote the action-less trajectory by $\tilde{\tau}_{:T}$, that is,

$$\tilde{\tau}_{:T} = (s_0, r_0, s_1, r_1, s_2, r_2, \dots, s_{T-1}, r_{T-1}, s_T).$$

By Jensen's inequality,

 $\log p(\tau_{:T})$

$$= \log \rho_{0}(s_{0}) + \log \sum_{m=0}^{K-1} p(m) \prod_{t=0}^{T-1} [P(s_{t+1} \mid s_{t}; m)\beta(a_{t} \mid s_{t}; m)R(r_{t} \mid s_{t}, s_{t+1}; m)]$$

$$= \log \rho_{0}(s_{0}) + \log \sum_{m=0}^{K-1} q(m) \cdot \frac{p(m)}{q(m)} \prod_{t=0}^{T-1} [P(s_{t+1} \mid s_{t}; m)\beta(a_{t} \mid s_{t}; m)R(r_{t} \mid s_{t}, s_{t+1}; m)]$$

$$\geq \log \rho_{0}(s_{0}) + \mathbb{E}_{m \sim q} \left[\log \frac{p(m)}{q(m)} + \sum_{t=0}^{T-1} A(s_{t}, a_{t}, r_{t}, s_{t+1}; m) \right]$$

$$= \log \rho_{0}(s_{0}) - D_{\text{KL}}(q \parallel p) + \sum_{t=0}^{T-1} \mathbb{E}_{m \sim q} [A(s_{t}, a_{t}, r_{t}, s_{t+1}; m)],$$

where

$$A(s_{t}, a_{t}, r_{t}, s_{t+1}; m) = \log P(s_{t+1} \mid s_{t}; m) + \log \beta(a_{t} \mid s_{t}; m) + \log R(r_{t} \mid s_{t}, s_{t+1}; m).$$

We restate the two assumptions we made in Section 4.2 for the reader's convenience.

Assumption 3. There is $K_P > 0$ such that for all $s \in \mathcal{S}$ and $a_1, a_2 \in \mathcal{A}$, $W_1(P(s, a_1), P(s, a_2)) < 0$ $K_P ||a_1 - a_2||$.

Assumption 4. The value function of the behaviour policy β is K_V -Lipschitz.

Lemma 6. For any policy π and $s \in \mathcal{S}$,

$$\frac{r_{\min}}{1 - \gamma} \le V^{\beta}(s) \le \frac{r_{\max}}{1 - \gamma}$$

Proof. By the definition of V^{π} , for all $s \in \mathcal{S}$,

$$V^{\pi}(s) = \mathbb{E}_{\tau \sim \pi|s} \left[\sum_{t=0}^{\infty} \gamma^{t} r_{t} \right] \ge \mathbb{E}_{\tau \sim \pi|s} \left[\sum_{t=0}^{\infty} \gamma^{t} r_{\min} \right] = \frac{r_{\min}}{1 - \gamma},$$

$$V^{\pi}(s) = \mathbb{E}_{\tau \sim \pi \mid s} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \leq \mathbb{E}_{\tau \sim \pi \mid s} \left[\sum_{t=0}^{\infty} \gamma^t r_{\max} \right] = \frac{r_{\max}}{1 - \gamma}.$$

Proposition 7. Define $Q^{\mathrm{LB}}_{\beta} \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ by the equation

$$Q_{\beta}^{\mathrm{LB}}(s,a) = \max \left\{ V^{\beta}(s) - r_{\max} + r_{\min} - \gamma K_V K_P \mathbb{E}_{a' \sim \beta(s)} \left[\|a - a'\| \right], \frac{r_{\min}}{1 - \gamma} \right\}.$$

For any policy $\pi \colon \mathcal{S} \to \mathcal{P}(\mathcal{A})$ such that $V^{\pi} \succeq V^{\beta}$, $Q^{\pi} \succeq Q^{\mathrm{LB}}_{\beta}$.

Proof. By Lemma 6 and the definition of Q^{π} , for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q^{\pi}(s, a) = \mathbb{E}_{s' \sim P(s, a, s'), r \sim R(s, a, s')} \left[r + \gamma V^{\pi}(s') \right]$$

$$\geq \mathbb{E}_{s' \sim P(s, a, s'), r \sim R(s, a, s')} \left[r_{\min} + \gamma \frac{r_{\min}}{1 - \gamma} \right]$$

$$= \frac{r_{\min}}{1 - \gamma}.$$

So we only need to show that for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q^{\pi}(s, a) \ge V^{\beta}(s) - r_{\max} + r_{\min} - \gamma K_V K_P \mathbb{E}_{a' \sim \beta(s)} [\|a - a'\|].$$

Let $a_1, a_2 \in \mathcal{A}$. By the Kantorovich–Rubinstein formula (Villani, 2009),

$$\left| \mathbb{E}_{s' \sim P(\cdot \mid s, a_1)} [V^{\beta}(s')] - \mathbb{E}_{s' \sim P(\cdot \mid s, a_2)} [V^{\beta}(s')] \right| \le K_V W_1(P(\cdot \mid s, a_1), P(\cdot \mid s, a_2))$$

$$\le K_V K_P \|a_1 - a_2\|.$$

Therefore,

$$Q^{\pi}(s,a) = \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} [r + \gamma V^{\pi}(s')]$$

$$\geq \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} [r + \gamma V^{\beta}(s')]$$

$$\geq V^{\beta}(s) - \mathbb{E}_{a' \sim \beta(\cdot|s)} [\mathbb{E}_{s' \sim P(s,a'),r \sim R(s,a',s')} [r + \gamma V^{\beta}(s')]]$$

$$+ \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} [r + \gamma V^{\beta}(s')]$$

$$\geq V^{\beta}(s) - r_{\max} + r_{\min} + \gamma \mathbb{E}_{a' \sim \beta(s)} [\mathbb{E}_{s' \sim P(s,a)} [V^{\beta}(s')] - \mathbb{E}_{s' \sim P(s,a')} [V^{\beta}(s')]]$$

$$\geq V^{\beta}(s) - r_{\max} + r_{\min} - \gamma K_{V} K_{P} \mathbb{E}_{a' \sim \beta(\cdot|s)} [||a - a'||].$$

Note that we have used the fact that

$$V^{\beta}(s) = \mathbb{E}_{a' \sim \beta(s), s' \sim P(s, a'), r \sim R(s, a', s')} \left[r + \gamma V^{\beta}(s') \right].$$

Theorem 8. Let $\{A_s\}_{s\in\mathcal{S}}$ be a family of subsets of \mathcal{A} , $Q\in L^{\infty}(\mathcal{S}\times\mathcal{A})$, and \mathcal{T}_A be an operator on the space of real-valued functions on $\mathcal{S}\times\mathcal{A}$ defined by the equation

$$(\mathcal{T}_A Q)(s,a) = \begin{cases} (\mathcal{T}^* Q)(s,a) & \text{if } a \in A_s, \\ \tilde{Q}(s,a) & \text{otherwise,} \end{cases}$$

for each $Q \in L^{\infty}(S \times A)$. Then any bounded real-valued function on $S \times A$ converges to a unique fixed point Q_A by repeatedly applying \mathcal{T}_A .

Proof. Fix $s \in \mathcal{S}$, $a \in \mathcal{A}$, and $Q \in L^{\infty}(\mathcal{S} \times \mathcal{A})$. If $a \in A_s$,

$$\begin{aligned} |(\mathcal{T}_{A}Q)(s,a)| &= |(\mathcal{T}^{*}Q)(s,a)| \\ &= \left| \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} \left[r + \gamma \sup_{a' \in \mathcal{A}} Q(s',a') \right] \right| \\ &\leq \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} \left[\left| r + \gamma \sup_{a' \in \mathcal{A}} Q(s',a') \right| \right] \\ &\leq \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} \left[\left| r \right| + \gamma \left| \sup_{a' \in \mathcal{A}} Q(s',a') \right| \right] \\ &\leq \mathbb{E}_{s' \sim P(s,a),r \sim R(s,a,s')} \left[\max \left\{ \left| r_{\max} \right|, \left| r_{\min} \right| \right\} + \gamma \|Q\|_{\infty} \right], \\ &= \max \left\{ \left| r_{\max} \right|, \left| r_{\min} \right| \right\} + \gamma \|Q\|_{\infty}. \end{aligned}$$

Otherwise,

$$|(\mathcal{T}_A Q)(s,a)| = |\tilde{Q}(s,a)| \le ||\tilde{Q}||_{\infty}.$$

So

$$\|\mathcal{T}_A Q\|_{\infty} \le \max\left\{\|\tilde{Q}\|_{\infty}, \max\left\{|r_{\max}|, |r_{\min}|\right\} + \gamma\|Q\|_{\infty}\right\} < \infty,$$

that is, $\mathcal{T}_A Q \in L^{\infty}(\mathcal{S} \times \mathcal{A})$. So the restriction of \mathcal{T}_A onto $L^{\infty}(\mathcal{S} \times \mathcal{A})$ is an operator on $L^{\infty}(\mathcal{S} \times \mathcal{A})$. With a slight abuse of notation, we will just denote the restriction by \mathcal{T}_A from now on.

Now we go on and prove that \mathcal{T}_A is a contraction operator. Fix $s \in \mathcal{S}$, $a \in \mathcal{A}$ and $Q_1, Q_2 \in L^{\infty}(\mathcal{S} \times \mathcal{A})$. If $a \in A_s$,

$$|(\mathcal{T}_{A}Q_{1})(s,a) - (\mathcal{T}_{A}Q_{2})(s,a)| = |(\mathcal{T}^{*}Q_{1})(s,a) - (\mathcal{T}^{*}Q_{2})(s,a)|$$

$$= \gamma \left| \mathbb{E}_{s' \sim P(s,a)} \left[\sup_{a' \in \mathcal{A}} Q_{1}(s',a') - \sup_{a'' \in \mathcal{A}} Q_{2}(s',a'') \right] \right|$$

$$\leq \gamma \mathbb{E}_{s' \sim P(s,a)} \left[\left| \sup_{a' \in \mathcal{A}} Q_{1}(s',a') - \sup_{a'' \in \mathcal{A}} Q_{2}(s',a'') \right| \right]$$

$$\leq \gamma \mathbb{E}_{s' \sim P(s,a)} \left[\sup_{a' \in \mathcal{A}} |Q_{1}(s',a') - Q_{2}(s',a')| \right]$$

$$\leq \gamma \|Q_{1} - Q_{2}\|_{\infty}.$$

Otherwise,

$$|(\mathcal{T}_A Q_1)(s, a) - (\mathcal{T}_A Q_2)(s, a)| = |\tilde{Q}(s, a) - \tilde{Q}(s, a)| = 0 \le \gamma ||Q_1 - Q_2||_{\infty}.$$

Therefore, and \mathcal{T}_A is a contraction mapping on $L^{\infty}(\mathcal{S} \times \mathcal{A})$. By the contraction mapping theorem, any initial-bounded Q-function would converge to a unique fixed point Q_A .

Lemma 9. Let π_1 and π_2 be two policies. If $\mathbb{E}_{a \sim \pi_1(s)}[Q^{\pi_2}(s, a)] \geq V^{\pi_2}(s)$ for all $s \in \mathcal{S}$, then $V^{\pi_1} \succeq V^{\pi_2}$.

Proof. We define a sequence (Q_n) of bounded real-valued functions on $\mathcal{S} \times \mathcal{A}$ by the recurrence relation

$$Q_n = \begin{cases} Q^{\pi_2} & \text{if } n = 0, \\ \mathcal{T}^{\pi_1} Q_{n-1} & \text{otherwise.} \end{cases}$$

We first show that $Q_n \succeq Q^{\pi_2}$ by mathematical induction. The base case is trivial because $Q_0 \equiv Q^{\pi_2}$. Suppose $Q_{n-1} \succeq Q^{\pi_2}$. Then for each $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q_{n}(s, a) = (\mathcal{T}^{\pi_{1}}Q_{n-1})(s, a)$$

$$= \mathbb{E}_{s' \sim P(s, a), r \sim R(s, a, s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi_{1}(s')} \left[Q_{n-1}(s', a') \right] \right]$$

$$\geq \mathbb{E}_{s' \sim P(s, a), r \sim R(s, a, s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi_{1}(s')} \left[Q^{\pi_{2}}(s', a') \right] \right]$$

$$\geq \mathbb{E}_{s' \sim P(s, a), r \sim R(s, a, s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi_{2}(s')} \left[Q^{\pi_{2}}(s', a') \right] \right]$$

$$= (\mathcal{T}^{\pi_{2}}Q^{\pi_{2}})(s, a)$$

$$= Q^{\pi_{2}}(s, a).$$

So $Q_n \succeq Q^{\pi_2}$. By mathematical induction, $Q_n \succeq Q^{\pi_2}$ for all n. For all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q^{\pi_1}(s, a) = \lim_{n \to \infty} Q_n(s, a) \ge Q^{\pi_2}(s, a)$$

Therefore, for all $s \in \mathcal{S}$,

$$V^{\pi_1}(s) = \mathbb{E}_{a \sim \pi_1(s)} \left[Q^{\pi_1}(s,a) \right] \ge \mathbb{E}_{a \sim \pi_1(s)} \left[Q^{\pi_2}(s,a) \right] \ge V^{\pi_2}(s),$$
 that is, $V^{\pi_1} \succ V^{\pi_2}$.

Theorem 10. Let $\{A_s\}_{s\in\mathcal{S}}$ be a family of subsets of \mathcal{A} , $\tilde{Q}\in L^{\infty}(\mathcal{S}\times\mathcal{A})$, and Q_A be a bounded real-valued function that satisfies the relation

$$Q_A(s,a) = \begin{cases} (\mathcal{T}^*Q_A)(s,a) & \text{if } a \in A_s, \\ \tilde{Q}(s,a) & \text{otherwise,} \end{cases}$$
 (10)

for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$. Suppose there is a policy π such that for all $s \in \mathcal{S}$,

$$V^{\pi}(s) \ge \sup_{a \in \mathcal{A}} \tilde{Q}(s, a)$$

and

$$Q^{\pi}(s, a_s) \ge V^{\pi}(s)$$

for some $a_s \in A_s$. If there exists a deterministic policy $\pi_A^* : S \to A$ that is optimal under the constraint $\pi_A(s) \in A_s$ for all $s \in S$, then

$$\pi_A^*(s) = \arg\max_{a \in A} Q_A(s, a).$$

Proof. Define \mathcal{T}_A as in Theorem 8. We can see that there is a unique bounded real-valued function Q_A that satisfies (10), because by Theorem 8, \mathcal{T}_A has unique fixed point Q_A .

We proceed to prove that for each $s \in \mathcal{S}$, $Q_A(s, \pi_A^*(s)) \geq V^{\pi_A^*}(s)$. Define a sequence (Q_n) of bounded real-valued functions on $\mathcal{S} \times \mathcal{A}$ by the recurrence relation

$$Q_n = \begin{cases} Q_0 & \text{if } n = 0, \\ \mathcal{T}_A Q_{n-1} & \text{otherwise,} \end{cases}$$
 (11)

where $Q_0 : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is defined as

$$Q_0(s,a) = \begin{cases} Q^{\pi_A^*}(s,a) & \text{if } a \in A_s, \\ \tilde{Q}(s,a) & \text{otherwise}. \end{cases}$$

When n = 0, for all $s \in \mathcal{S}$

$$Q_0(s, \pi_A^*(s)) = Q^{\pi_A^*}(s, \pi_A^*(s)) = V^{\pi_A^*}(s),$$

because $\pi_A^*(s) \in A_s$. Assume $Q_{n-1}(s, \pi_A^*(s)) \geq V^{\pi_A^*}(s)$ for all $s \in \mathcal{S}$. Then for all $s \in \mathcal{S}$,

$$Q_{n}(s, \pi_{A}^{*}(s)) = \mathbb{E}_{s' \sim P(s, \pi_{A}^{*}(s)), r \sim R(s, \pi_{A}^{*}(s), s')} \left[r + \gamma \sup_{a' \in A} Q_{n-1}(s', a') \right]$$

$$\geq \mathbb{E}_{s' \sim P(s, \pi_{A}^{*}(s)), r \sim R(s, \pi_{A}^{*}(s), s')} \left[r + \gamma Q_{n-1}(s', \pi_{A}^{*}(s')) \right]$$

$$\geq \mathbb{E}_{s' \sim P(s, \pi_{A}^{*}(s)), r \sim R(s, \pi_{A}^{*}(s), s')} \left[r + \gamma V^{\pi_{A}^{*}}(s') \right]$$

$$= \mathbb{E}_{s' \sim P(s, \pi_{A}^{*}(s)), r \sim R(s, \pi_{A}^{*}(s), s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi_{A}^{*}(s')} \left[Q^{\pi_{A}^{*}}(s', a') \right] \right]$$

$$= (\mathcal{T}^{\pi_{A}^{*}} Q^{\pi_{A}^{*}})(s, \pi_{A}^{*}(s))$$

$$= Q^{\pi_{A}^{*}}(s, \pi_{A}^{*}(s))$$

$$= V^{\pi_{A}^{*}}(s).$$

So by mathematical induction, $Q_n(s, \pi_A^*(s)) \geq V^{\pi_A^*}(s)$ for all $s \in \mathcal{S}$ and $n \geq 0$. Therefore,

$$Q_A(s, \pi_A^*(s)) = \lim_{n \to \infty} Q_n(s, \pi_A^*(s)) \ge V^{\pi_A^*}(s).$$

Since for all $s \in \mathcal{S}$ and $a \in A_s$,

$$Q_A(s,a) = (\mathcal{T}_A Q_A)(s,a) = \tilde{Q}(s,a) \le V^{\pi}(s).$$

We can define a deterministic policy $\pi_A \colon \mathcal{S} \to \mathcal{A}$ that maps $s \in \mathcal{S}$ to a_s . Since $\pi_A(s) = a_s \in A_s$ for all $s \in \mathcal{S}$ and π_A^* is the optimal policy among the policies that satisfy this constraint, we have $V^{\pi_A^*} \succeq V^{\pi_A}$. So we may conclude that for all $s \in \mathcal{S}$,

$$\sup_{a \in \mathcal{A} \setminus A_s} Q_A(s, a) \le V^{\pi}(s) \le V^{\pi_A}(s) \le V^{\pi_A^*}(s) = Q_A(s, \pi_A^*(s)). \tag{12}$$

We finish the proof by showing that for all $s \in \mathcal{S}$,

$$Q_A(s, \pi_A^*(s)) = \max_{a \in A} Q_A(s, a).$$

Recall the sequence (Q_n) we previously defined by the recurrence relation (11). We will prove that for every $n, s \in \mathcal{S}$, and $a \in \mathcal{A}$,

$$Q_n(s,a) \le V^{\pi_A^*}(s) = Q_n(s, \pi_A^*(s)).$$

Assume n=0. Fix $s^{\dagger} \in \mathcal{S}$ and $a^{\dagger} \in \mathcal{A}$. If $a^{\dagger} \notin A_{s^{\dagger}}$, then by the observation we made in (12),

$$Q_0(s^{\dagger}, a^{\dagger}) = \tilde{Q}(s^{\dagger}, a^{\dagger}) \le V^{\pi_A^*}(s^{\dagger}).$$

If $a^{\dagger} \in A_{s^{\dagger}}$, consider a policy $\pi^{\dagger} \colon \mathcal{S} \to \mathcal{A}$ defined as

$$\pi^{\dagger}(s) = \begin{cases} a^{\dagger} & \text{if } s = s^{\dagger}, \\ \pi_{A}^{*}(s) & \text{otherwise.} \end{cases}$$

For all $s \in \mathcal{S}$, $\pi^{\dagger}(s) \in A_s$, so $V^{\pi_A^*} \succeq V^{\pi^{\dagger}}$. If $Q_0(s^{\dagger}, a^{\dagger}) < V^{\pi_A^*}(s^{\dagger})$, it satisfies our hypothesis. Otherwise.

$$Q^{\pi_A^*}(s^{\dagger}, \pi^{\dagger}(s^{\dagger})) = Q_0(s^{\dagger}, a^{\dagger}) \ge V^{\pi_A^*}(s^{\dagger}) = Q^{\pi_A^*}(s^{\dagger}, \pi_A^*(s^{\dagger}))$$

and for $s \neq s^{\dagger}$,

$$Q^{\pi_A^*}(s, \pi^{\dagger}(s)) = Q^{\pi_A^*}(s, \pi_A^*(s)),$$

so by Lemma 9, $V^{\pi^{\dagger}} \succeq V^{\pi_A^*}$, which means $V^{\pi^{\dagger}} \equiv V^{\pi_A^*}$. Then

$$\begin{split} Q^{\pi_A^*}(s^\dagger, a^\dagger) &= (\mathcal{T}^{\pi_A^*} Q^{\pi_\beta^*})(s^\dagger, a^\dagger) \\ &= \mathbb{E}_{s' \sim P(s^\dagger, a^\dagger), r \sim R(s^\dagger, a^\dagger, s')} \left[r + \gamma Q^{\pi_A^*}(s', \pi_A^*(s')) \right] \\ &= \mathbb{E}_{s' \sim P(s^\dagger, a^\dagger), r \sim R(s^\dagger, a^\dagger, s')} \left[r + \gamma V^{\pi_A^*}(s') \right] \\ &= \mathbb{E}_{s' \sim P(s^\dagger, a^\dagger), r \sim R(s^\dagger, a^\dagger, s')} \left[r + \gamma V^{\pi^\dagger}(s') \right] \\ &= \mathbb{E}_{s' \sim P(s^\dagger, a^\dagger), r \sim R(s^\dagger, a^\dagger, s')} \left[r + \gamma Q^{\pi^\dagger}(s', \pi^\dagger(s')) \right] \\ &= (\mathcal{T}^{\pi^\dagger} Q^{\pi^\dagger})(s^\dagger, a^\dagger) \\ &= Q^{\pi^\dagger}(s^\dagger, \pi^\dagger(s^\dagger)) \\ &= V^{\pi_A^*}(s^\dagger). \end{split}$$

So $Q_0(s^{\dagger}, a^{\dagger}) \leq V^{\pi_A^*}(s^{\dagger})$ in both cases. Since it is obvious that

$$Q_0(s^{\dagger}, \pi_A^*(s^{\dagger})) = Q^{\pi_A^*}(s^{\dagger}, \pi_A^*(s^{\dagger})) = V^{\pi_A^*}(s^{\dagger}),$$

our hypothesis holds for n = 0.

Assume the hypothesis holds for n-1. Fix $s^{\dagger} \in \mathcal{S}$ and $a^{\dagger} \in \mathcal{A}$. If $a^{\dagger} \notin A_{s^{\dagger}}$,

$$Q_n(s^\dagger,a^\dagger) = (\mathcal{T}_\beta^*Q_{n-1})(s^\dagger,a^\dagger) = Q_\beta^{\mathrm{LB}}(s^\dagger,a^\dagger) \leq V^{\pi_\beta^*}(s^\dagger)$$

by (12). Otherwise,

$$\begin{aligned} Q_{n}(s^{\dagger}, a^{\dagger}) &= (\mathcal{T}_{A}Q_{n-1})(s^{\dagger}, a^{\dagger}) \\ &= (\mathcal{T}^{*}Q_{n-1})(s^{\dagger}, a^{\dagger}) \\ &= \mathbb{E}_{s' \sim P(s^{\dagger}, a^{\dagger}), r \sim R(s^{\dagger}, a^{\dagger}, s')} \left[r + \gamma \sup_{a' \in \mathcal{A}} Q_{n-1}(s', a') \right] \\ &= \mathbb{E}_{s' \sim P(s^{\dagger}, a^{\dagger}), r \sim R(s^{\dagger}, s^{\dagger}, s')} \left[r + \gamma \max_{a' \in \mathcal{A}} Q_{n-1}(s', a') \right] \\ &= \mathbb{E}_{s' \sim P(s^{\dagger}, a^{\dagger}), r \sim R(s^{\dagger}, s^{\dagger}, s')} \left[r + \gamma V^{\pi_{A}^{*}}(s') \right] \\ &= \mathbb{E}_{s' \sim P(s^{\dagger}, a^{\dagger}), r \sim R(s^{\dagger}, s^{\dagger}, s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi_{A}^{*}(s')} \left[Q^{\pi_{A}^{*}}(s', a') \right] \right] \\ &= (\mathcal{T}^{\pi_{A}^{*}} Q^{\pi_{A}^{*}})(s^{\dagger}, a^{\dagger}) \\ &= Q_{0}(s^{\dagger}, a^{\dagger}) \\ &= Q_{0}(s^{\dagger}, a^{\dagger}) \\ &< V^{\pi_{A}^{*}}(s^{\dagger}). \end{aligned}$$

When $a^{\dagger} = \pi_A^*(s^{\dagger})$, the inequality becomes equality. So by mathematical induction, for every n, $s \in \mathcal{S}$, and $a \in \mathcal{A}$,

$$Q_n(s, a) \le V^{\pi_A^*}(s) = Q_n(s, \pi_A^*(s)).$$

Sending n to infinity, we can see that for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q_A(s,a) = \lim_{n \to \infty} Q_n(s,a) \le V^{\pi_A^*}(s) = \lim_{n \to \infty} Q_n(s,\pi_A^*(s)) = Q_A(s,\pi_A^*(s)).$$

969 Therefore,

$$Q_A(s, \pi_A^*(s)) = \max_{a \in \mathcal{A}} Q_A(s, a).$$

Theorem 11. Any initial bounded real-valued function on $S \times A$ can converge to a unique fixed point Q_{β}^* by repeatedly applying \mathcal{T}_{β}^* . Suppose for each $s \in S$,

$$Q^{\beta}(s, a_s) \ge \mathbb{E}_{a \sim \beta(s)} \left[Q^{\beta}(s, a) \right]$$

for some $a_s \in \mathcal{A} \setminus \mathrm{OOD}(s)$. If there exists a deterministic policy $\pi_{\beta}^* \colon \mathcal{S} \to \mathcal{A}$ that is optimal under the constraint $\pi(s) \notin \mathrm{OOD}(s)$ for all $s \in \mathcal{S}$, then $\pi_{\beta}^*(s) = \arg\max_{a \in \mathcal{A}} Q_{\beta}^*(s, a)$ for all $s \in \mathcal{S}$.

Proof. First observe that for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q_{\beta}^{\mathrm{LB}}(s,a) \ge \frac{r_{\min}}{1-\gamma}$$

and

$$Q_{\beta}^{\text{LB}}(s, a) \le V^{\beta}(s) - r_{\text{max}} + r_{\text{min}} \le \frac{r_{\text{max}}}{1 - \gamma} - r_{\text{max}} + r_{\text{min}},$$

by Lemma 6. This implies $Q_{\beta}^{\mathrm{LB}} \in L^{\infty}(\mathcal{S} \times \mathcal{A})$. Since $\tilde{Q} = Q_{\beta}^{\mathrm{LB}}$ and $A_s = \mathcal{A} \setminus \mathrm{OOD}(s)$ satisfies the conditions of Lemma 8, any initially bounded real-valued function on $\mathcal{S} \times \mathcal{A}$ converges to a unique fixed point, which we denote by Q_{β}^* , through repeated application of \mathcal{T}_A , which is in fact, \mathcal{T}_{β}^* .

For all $s \in \mathcal{S}$,

$$Q_{\beta}^{\text{LB}}(s, a) \le V^{\beta}(s) - r_{\text{max}} + r_{\text{min}} \le V^{\beta}(s),$$

which means

$$\sup_{a \in \mathcal{A}} Q_{\beta}^{\mathrm{LB}}(s, a). \le V^{\beta}(s) = \mathbb{E}_{a \sim \beta(s)} \left[Q^{\beta}(s, a) \right] < \sup_{a \in \mathcal{A} \setminus \mathrm{OOD}(s)} Q^{\beta}(s, a).$$

Now we can see that the second part of the theorem is a special case of Theorem 10, where $A_s = \mathcal{A} \setminus \mathrm{OOD}(s)$, $\tilde{Q} = Q_{\beta}^{\mathrm{LB}}$, $Q_A = Q_{\beta}^*$, $\pi = \beta$, and $\pi_A^* = \pi_{\beta}^*$.

Lemma 12. Let $\Pi = \{\pi_0, \pi_1, \pi_2, \dots, \pi_{N-1}\}$ be a finite set of policies. If π^* is a policy such that for each $s \in S$, there is $i \in [N]$ such that $\pi^*(s) = \pi_i(s)$ and $V^{\pi_i}(s) = \max_{\pi \in \Pi} V^{\pi}(s)$, then $V^{\pi^*} \succeq V^{\pi}$ for every $\pi \in \Pi$.

Proof. Define a sequence (Q_n) of bounded real-valued functions by the recurrence relation

$$Q_n = \begin{cases} \max_{\pi \in \Pi} Q^{\pi} & \text{if } n = 0, \\ \mathcal{T}^{\pi^*} Q_{n-1}, & \text{otherwise.} \end{cases}$$

We want to show that $Q_n \succeq \max_{\pi \in \Pi} Q^{\pi}$ for all $n \geq 0$. The base case is trivial. Assume $Q_{n-1} \succeq \max_{\pi \in \Pi} Q^{\pi}$. For each $s \in \mathcal{S}$, there is $i \in [N]$ such that $\pi^*(s) = \pi_i(s)$ and $V^{\pi_i}(s) = \max_{\pi \in \Pi} V^{\pi}(s)$, which implies

$$\mathbb{E}_{a \sim \pi^*(s)} \left[Q_{n-1}(s, a) \right] \ge \mathbb{E}_{a \sim \pi_i(s)} \left[Q^{\pi_i}(s) \right] = V^{\pi_i}(s) = \max_{\pi \in \Pi} V^{\pi}(s).$$

Now for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$Q_{n}(s, a) = (\mathcal{T}^{\pi^{*}}Q_{n-1})(s, a)$$

$$= \mathbb{E}_{s' \sim P(s, a), r \sim R(s, a, s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi^{*}(s')} \left[Q_{n-1}(s', a') \right] \right]$$

$$\geq \mathbb{E}_{s' \sim P(s, a), r \sim R(s, a, s')} \left[r + \gamma \max_{\pi \in \Pi} V^{\pi}(s') \right]$$

$$= \max_{\pi \in \Pi} \mathbb{E}_{s' \sim P(s, a), r \sim R(s, a, s')} \left[r + \gamma \mathbb{E}_{a' \sim \pi(s')} \left[Q^{\pi}(s', a') \right] \right]$$

$$= \max_{\pi \in \Pi} (\mathcal{T}^{\pi}Q^{\pi})(s, a)$$

$$= \max_{\pi \in \Pi} Q^{\pi}(s, a).$$

By mathematical induction, $Q_n \succeq \max_{\pi \in \Pi} Q^{\pi}$ for all $n \geq 0$. Therefore,

$$Q^{\pi^*}(s, a) = \lim_{n \to \infty} Q_n(s, a) \ge \max_{\pi \in \Pi} Q^{\pi}(s, a)$$

for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$.

Fix $s \in \mathcal{S}$. There is $i \in [N]$ such that $\pi^*(s) = \pi_i(s)$ and $V^{\pi_i}(s) = \max_{\pi \in \Pi} V^{\pi}(s)$. Then

$$V^{\pi^*}(s) = \mathbb{E}_{a \sim \pi^*(s)} \left[Q^{\pi^*}(s, a) \right] \ge \mathbb{E}_{a \sim \pi_i(s)} \left[Q^{\pi_i}(s, a) \right] = V^{\pi_i}(s) = \max_{\pi \in \Pi} V^{\pi}(s).$$

Our choice of s was arbitrary, so $V^{\pi^*} \succeq \max_{\pi \in \Pi} V^{\pi}$.

Theorem 13. Any initial bounded real-valued function on $S \times A$ can converge to a unique fixed point Q_B^* by repeatedly applying \mathcal{T}_B^* . Suppose for each $\beta \in \mathcal{B}$ and $s \in S$,

$$Q^{\beta}(s, a_s^{\beta}) \ge \mathbb{E}_{a \sim \beta(s)} \left[Q^{\beta}(s, a) \right]$$

for some $a_s^{\beta} \in \mathcal{A} \setminus \mathrm{OOD}(s)$. If there exists a deterministic policy $\pi_{\mathcal{B}}^* \colon \mathcal{S} \to \mathcal{A}$ that is optimal under the constraint $\pi(s) \notin \mathrm{OOD}(s)$ for all $s \in \mathcal{S}$, then $\pi_{\mathcal{B}}^*(s) = \arg \max_{a \in \mathcal{A}} Q_{\mathcal{B}}^*(s, a)$ for all $s \in \mathcal{S}$.

Proof. First observe that for all $s \in \mathcal{S}$, $a \in \mathcal{A}$, and $\beta \in \mathcal{B}$,

$$Q_{\beta}^{\mathrm{LB}}(s, a) \ge \frac{r_{\min}}{1 - \gamma}$$

and

$$Q_{\beta}^{\mathrm{LB}}(s,a) \leq V^{\beta}(s) - r_{\mathrm{max}} + r_{\mathrm{min}} \leq \frac{r_{\mathrm{max}}}{1 - \gamma} - r_{\mathrm{max}} + r_{\mathrm{min}},$$

by Lemma 6. So obviously,

$$Q_{\mathcal{B}}^{\mathrm{LB}}(s, a) = \max_{\beta \in \mathcal{B}} Q_{\beta}^{\mathrm{LB}}(s, a) \ge \frac{r_{\min}}{1 - \gamma},$$

and

$$Q_{\mathcal{B}}^{\text{LB}}(s, a) = \max_{\beta \in \mathcal{B}} Q_{\beta}^{\text{LB}}(s, a) \le \frac{r_{\text{max}}}{1 - \gamma} - r_{\text{max}} + r_{\text{min}},$$

for all $s \in \mathcal{S}$ and \mathcal{A} . This implies $Q_{\mathcal{B}}^{\mathrm{LB}} \in L^{\infty}(\mathcal{S} \times \mathcal{A})$. Since $\tilde{Q} = Q_{\mathcal{B}}^{\mathrm{LB}}$ and $A_s = \mathcal{A} \setminus \mathrm{OOD}(s)$ satisfies the conditions of Lemma 8, any initially bounded real-valued function on $\mathcal{S} \times \mathcal{A}$ converges to a unique fixed point, which we denote by $Q_{\mathcal{B}}^*$, through repeated application of \mathcal{T}_A , which is in fact, $\mathcal{T}_{\mathcal{B}}^*$.

For each $\beta \in \mathcal{B}$ define a deterministic policy $\beta' \colon \mathcal{S} \to \mathcal{A}$ so that $\beta'(s) = a_s^{\beta}$ for each $s \in \mathcal{S}$. Then

$$Q^{\beta}(s, \beta'(s)) = Q^{\beta}(s, a_s^{\beta}) \ge \mathbb{E}_{a \sim \beta(s)} \left[Q^{\beta}(s, a) \right]$$

for all $s \in \mathcal{S}$, so $V^{\beta'} \succeq V^{\beta}$ by Lemma 9. We will denote the set $\{\beta' : \beta \in \mathcal{B}\}$ by \mathcal{B}' . Consider a policy $\beta^* : \mathcal{S} \to \mathcal{A}$ defined as

$$\beta^*(s) = \left(\underset{\beta' \in \mathcal{B}'}{\arg \max} V^{\beta'}(s)\right)(s),$$

that is, for each state, we follow the β' with the highest value. Obviously, $\beta^*(s) \in \mathcal{A} \setminus \text{OOD}(s)$ for all $s \in \mathcal{S}$, and by Lemma 12, for all $s \in \mathcal{S}$, $a \in \mathcal{A}$, and $\beta \in \mathcal{B}$,

$$Q_{\beta}^{\mathrm{LB}}(s,a) \leq V^{\beta}(s) - r_{\mathrm{max}} + r_{\mathrm{min}} \leq V^{\beta}(s) \leq V^{\beta'}(s) \leq V^{\beta^*}(s),$$

which means

$$V^{\beta^*}(s) \ge \sup_{a \in A} \max_{\beta \in \mathcal{B}} Q^{\mathrm{LB}}_{\beta}(s, a) = \sup_{a \in A} Q^{\mathrm{LB}}_{\mathcal{B}}(s, a).$$

Now we can see that the second part of the theorem is a special case of Theorem 10, where $A_s = \mathcal{A} \setminus \mathrm{OOD}(s)$, $\tilde{Q} = Q_{\mathcal{B}}^{\mathrm{LB}}$, $Q_A = Q_{\mathcal{B}}^*$, $\pi = \beta^*$, and $\pi_A^* = \pi_\beta^*$.

Proposition 14. Let X be a non-degenerate multivariate Gaussian random vector with mean $\mu \in \mathbb{R}^d$ and a diagonal covariance matrix $\operatorname{diag}(\sigma)^2 \in \mathbb{R}^{d \times d}$. For $y \in \mathbb{R}^d$,

$$\mathbb{E}\left[\left\|\mathbf{X} - \mathbf{y}\right\|_{1}\right] = \sum_{i=1}^{d} \left[(y_{i} - \mu_{i}) \operatorname{erf}\left(\frac{y_{i} - \mu_{i}}{\sigma_{i} \sqrt{2}}\right) + \sqrt{\frac{2}{\pi}} \sigma_{i} \exp\left(-\frac{(y_{i} - \mu_{i})^{2}}{2\sigma_{i}^{2}}\right).\right]$$

Proof. Let $\mathbf{X}=(X_1,X_2,\ldots,X_d)$, $\mathbf{y}=(y_1,y_2,\ldots,y_d)$, $\boldsymbol{\mu}=(\mu_1,\mu_2,\ldots,\mu_d)$, and $\boldsymbol{\sigma}=(\sigma_1,\sigma_2,\ldots,\sigma_d)$. We may assume that $\sigma_1,\sigma_2,\ldots,\sigma_d>0$. Then

$$\mathbb{E}\left[\left\|\mathbf{X} - \mathbf{y}\right\|_{1}\right] = \mathbb{E}\left[\sum_{i=1}^{d}\left|X_{i} - y_{i}\right|\right] = \sum_{i=1}^{d}\mathbb{E}\left[\left|X_{i} - y_{i}\right|\right].$$

Define $g_i(y) = \mathbb{E}[|X_i - y|].$

$$g'_i(y) = \mathbb{E}\left[\frac{\mathrm{d}}{\mathrm{d}y}|X_i - y|\right] = \mathbb{E}[\mathbf{1}_{X_i < y} - \mathbf{1}_{X_i > y}] = F_{X_i}(y) - (1 - F_{X_i}(y)) = 2F_{X_i}(y) - 1,$$

where F_{X_i} is the cumulative distribution function of X_i . So

$$g_i'(y) = \operatorname{erf}\left(\frac{y - \mu_i}{\sigma_i \sqrt{2}}\right).$$

Observe that

$$\begin{split} g_i(\mu_i) &= \mathbb{E}\left[|X_i - \mu_i|\right] \\ &= \frac{1}{\sqrt{2\pi}\sigma_i} \int_{\mu_i}^{\infty} (x_i - \mu_i) \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \mathrm{d}x_i \\ &- \frac{1}{\sqrt{2\pi}\sigma_i} \int_{-\infty}^{\mu_i} (x_i - \mu_i) \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}\right) \mathrm{d}x_i. \end{split}$$

Substituting $u_i = (x_i - \mu_i)/\sigma_i$,

$$g_{i}(\mu_{i}) = \frac{1}{\sqrt{2\pi}\sigma_{i}} \left[\int_{0}^{\infty} \sigma_{i} u_{i} e^{-\frac{1}{2}u_{i}^{2}} \sigma_{i} \, du_{i} - \int_{-\infty}^{0} \sigma_{i} u_{i} e^{-\frac{1}{2}u_{i}^{2}} \sigma_{i} \, du_{i} \right]$$

$$= \frac{\sigma_{i}}{\sqrt{2\pi}} \left[\int_{0}^{\infty} u_{i} e^{-\frac{1}{2}u_{i}^{2}} \, du_{i} - \int_{-\infty}^{0} u_{i} e^{-\frac{1}{2}u_{i}^{2}} \, du_{i} \right]$$

$$= \sqrt{\frac{2}{\pi}} \sigma_{i}.$$

By the fundamental theorem of calculus,

$$g_i(y) = g_i(\mu_i) + \int_{\mu_i}^y g_i'(v) \, dv = \sqrt{\frac{2}{\pi}} \sigma_i + \int_{\mu_i}^y \operatorname{erf}\left(\frac{v - \mu_i}{\sigma_i \sqrt{2}}\right) \, dv.$$

Substituting $z = (v - \mu_i)/(\sigma_i \sqrt{2})$,

$$\begin{split} \int_{\mu_i}^y \operatorname{erf}\left(\frac{v-\mu_i}{\sigma_i\sqrt{2}}\right) \, \mathrm{d}v &= \int_0^{(y-\mu_i)/(\sigma_i\sqrt{2})} \operatorname{erf}(z)\sigma_i\sqrt{2} \, \, \mathrm{d}z \\ &= \sqrt{2}\sigma_i \left[\left(\frac{y-\mu_i}{\sigma_i\sqrt{2}}\right) \operatorname{erf}\left(\frac{y-\mu_i}{\sigma_i\sqrt{2}}\right) + \frac{1}{\sqrt{\pi}} \exp\left(-\frac{(y-\mu_i)^2}{2\sigma_i^2}\right) - \frac{1}{\sqrt{\pi}} \right]. \end{split}$$

Therefore,

$$g_i(y) = (y - \mu_i) \operatorname{erf}\left(\frac{y - \mu_i}{\sigma_i \sqrt{2}}\right) + \sqrt{\frac{2}{\pi}} \sigma_i \exp\left(-\frac{(y - \mu_i)^2}{2\sigma_i^2}\right),$$

which implies

$$\mathbb{E}\left[\left\|\mathbf{X} - \mathbf{y}\right\|_{1}\right] = \sum_{i=1}^{d} g_{i}(y_{i}) = \sum_{i=1}^{d} \left[(y_{i} - \mu_{i}) \operatorname{erf}\left(\frac{y_{i} - \mu_{i}}{\sigma_{i} \sqrt{2}}\right) + \sqrt{\frac{2}{\pi}} \sigma_{i} \exp\left(-\frac{(y_{i} - \mu_{i})^{2}}{2\sigma_{i}^{2}}\right).\right]$$

C PRACTICAL ALGORITHM

C.1 STAGE I: BEHAVIOUR POLICY LEARNING

In theory, each behaviour policy is well-defined on every state $s \in \mathcal{S}$. However, in practice, we can trust our estimations only in the vicinity of the states they were trained on. The problem is, we train each $\hat{\beta}$ only on the states they are assigned to. Therefore, we need a mechanism to determine which behaviour policy estimates we can trust given a state $s \in \mathcal{S}$. For this purpose, we additionally train a classifier $f_{\psi} \colon \mathcal{S} \to \mathcal{P}_d([K])$ using the computed assignments and determine the credible set by the equation

$$\hat{\mathcal{V}}(s) = \left\{ \hat{\beta}_i \in \hat{\mathcal{B}} : f_{\psi}(i \mid s) \ge \frac{1}{b} \max_{j \in [K]} f_{\psi}(j \mid s) \right\},\,$$

where b>0 is a hyperparameter. We accordingly modify the definition of $Q_{\hat{\mathcal{B}}}^{\mathrm{LB}}$ to

$$Q_{\hat{\mathcal{B}}}^{\mathrm{LB}}(s,a) = \max_{\hat{\beta} \in \hat{\mathcal{V}}(s)} Q_{\hat{\beta}}^{\mathrm{LB}}(s,a),$$

for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$.

C.2 STAGE II: BEHAVIOUR VALUE LEARNING

To learn the value functions of all K behaviour policies in parallel, we leverage a network $V_{\zeta} \colon \mathcal{S} \to \mathbb{R}^K$ with K outputs. The per sample temporal difference (TD) loss function can be written by the equation

$$\ell_V(\zeta; s, r, s') = \left(V_{\zeta}(s)[\mathbb{A}(s)] - r - \gamma V_{\zeta'}(s')[\mathbb{A}(s')]\right)^2,\tag{13}$$

where s, r, and s' are the state, reward, and next state sampled from the dataset, respectively, ζ' is the target network parameter that is updated by polyak averaging as in Lillicrap et al. (2016), and $V_{\zeta}(s)[i]$ is the i-th coordinate of $V_{\zeta}(s)$. Note that $\mathbb{A}(s)$, the cluster assignment of s, is equal to $\mathbb{A}(s')$, because we assign each trajectory to the same cluster.

C.3 STAGE III: POLICY LEARNING

In practice, it is infeasible to compute the superior term in the penalised Bellman operator $\mathcal{T}_{\mathcal{B}}^*$. We instead adopt the actor-critic formulation that alternates between the policy improvement step and the critic learning step. The goal of the policy improvement step is to find an action that maximises the critic for each state. The challenge is that the critic is highly non-convex due to the penalisation of critic values for actions between the means of the behaviour policies, causing gradient methods to yield suboptimal solutions. Hence, we search for the optimal action in the vicinity of each behaviour policy's mean simultaneously. This is done by training a network to output not the optimal action itself but the difference between the optimal action and one of the behaviour policy means. Using a network $g_{\psi_{\pi}}: \mathcal{S} \to \mathbb{R}^{2d_a}$, we encode a Gaussian distribution $\tilde{\pi}$ with a d_a -dimensional mean vector and $d_a \times d_a$ diagonal covariance matrix. For a given state s, we choose the best behaviour policy β^* among $\hat{\mathcal{V}}(s)$ with respect to the current critic function $Q_{\psi_{\Omega}}$, that is,

$$\hat{\beta}^* = \arg\max_{\hat{\beta} \in \hat{\mathcal{V}}(s)} Q_{\psi_Q}(s, \boldsymbol{\mu}_{\hat{\beta}}(s) + \delta(s)),$$

where $\mu_{\hat{\beta}}$ is the mean vector of the behaviour policy estimate $\hat{\beta}$ and δ is a vector sampled from $\tilde{\pi}$. Our current policy can be computed according to the following equation:

$$\pi(s) = \boldsymbol{\mu}_{\hat{\beta}*}(s) + \delta.$$

The loss function for the policy improvement step can be written by the following equation:

$$\ell_{\pi}(\psi_{\pi};s) = -\max_{\hat{\beta} \in \hat{\mathcal{V}}(s)} Q_{\psi_{Q}}\left(s, \boldsymbol{\mu}_{\hat{\beta}}(s) + g_{\psi_{\pi}}(s)\right).$$

The critic learning step has two objectives: minimising the TD error and penalising the OOD actions. For the first objective, we adopt the conventional TD loss adapted to match the way we defined our policy, which is represented by the equation

$$\ell_Q^{\text{TD}}(\psi_Q; s, a, r, s') = (Q_{\psi_Q}(s, a) - T(r, s'))^2,$$

where s, a, r, and s' are the state, action, reward, and next state sampled from the dataset, respectively, and the TD target T(r, s') is defined as

$$T(r,s') = r + \gamma \max_{\hat{\beta} \in \hat{\mathcal{V}}(s)} Q_{\psi_Q'} \left(s', \boldsymbol{\mu}_{\hat{\beta}}(s') + g_{\psi_{\pi}'}(s') \right).$$

 ψ_Q' and ψ_π' in the preceding equation are the target critic network parameters and target actor network parameters, respectively, which are updated by polyak averaging to gradually follow ψ_Q and ψ_π .

The second objective of the critic learning step is to penalise the critic values of OOD actions towards $Q_{\hat{\mathcal{B}}}^{\mathrm{LB}}(s,a)$. To achieve this goal, we need to be able to sample an action a from $\mathcal{U}(\mathrm{OOD}(s))$ and compute $Q_{\hat{\mathcal{B}}}^{\mathrm{LB}}(s,a)$. In general, it is difficult to sample uniformly from an arbitrary set, so we sample an action from \mathcal{A} instead and ignore it if it does not lie in $\mathrm{OOD}(s)$. Since we are working with actions mapped via the inverse hyperbolic tangent function, we need to sample from $\tanh^{-1}(\mathcal{A}) = \mathbb{R}^{d_a}$, which is unbounded. We circumvent this issue by sampling from a sufficiently large hypercube $[-L, L]^{d_a}$. In particular, we set L=20 based on the observation that $\tanh(10)\approx 1.00$ under the 32-bit floating point representation, where we doubled 10 to secure a safety margin.

In order to compute $Q_{\hat{\beta}}^{\mathrm{LB}}$, we need to compute $Q_{\hat{\beta}}^{\mathrm{LB}}$ for each $\hat{\beta} \in \hat{\mathcal{B}}$. However, computing $Q_{\hat{\beta}}^{\mathrm{LB}}$ is not straightforward, due to the term $\mathbb{E}_{a'\sim\hat{\beta}(s)}[\|a-a'\|]$ in (7). We discovered that if $\hat{\beta}$ has a diagonal covariance matrix and we use a 1-norm, the expectation has the following closed-form expression (Proposition 14):

$$\mathbb{E}_{a' \sim \hat{\beta}(s)} \left[\|a - a'\|_1 \right] = \sum_{i=1}^{d_a} (y_i - \mu_i(s)) \operatorname{erf} \left(\frac{a_i - \mu_i(s)}{\sigma_i(s)\sqrt{2}} \right) + \sqrt{\frac{2}{\pi}} \sum_{i=1}^{d_a} \sigma_i \exp \left(-\frac{(a_i - \mu_i(s))^2}{2\sigma_i^2(s)} \right),$$

where $a=(a_1,a_2,\ldots,a_{d_a})$ and $\beta(s)$ is a Gaussian distribution with a state-dependent mean vector $\mu(s)=(\mu_1(s),\mu_2(s),\ldots,\mu_{d_a}(s))$ and a state-dependent covariance matrix whose main diagonal is $\sigma(s)=(\sigma_1(s),\sigma_2(s),\ldots,\sigma_{d_a}(s))$. For r_{\min} and r_{\max} , following Mao et al. (2023), we estimate them by the minimum and maximum rewards in all of the datasets of a given task, that is, for example r_{\min} and r_{\max} for a hopper-v2 dataset is computed by the minimum and maximum of the rewards in hopper-expert-v2, hopper-medium-v2, and hopper-random-v2. To sum up, the loss function for regularisation is

$$\ell_Q^{\text{reg}}(\psi_Q; \tilde{s}, \tilde{a}) = \mathbf{1}_{\tilde{a} \in \text{OOD}(\tilde{s})} \left(Q_{\psi_Q}(\tilde{s}, \tilde{a}) - Q_{\hat{B}}^{\text{LB}}(\tilde{s}, \tilde{a}) \right)^2,$$

where \tilde{s} is a state sampled from the dataset and \tilde{a} is an action sampled from π_{alg} . The resulting total loss can be written by the following equation

$$\ell_Q(\psi_Q) = \ell_Q^{\text{TD}}(\psi_Q) + w_Q \ell_Q^{\text{reg}}(\psi_Q), \tag{14}$$

where w_Q is a tuneable hyperparameter. We have omitted the samples in the preceding equation for simplicity.

For π_{alg} we adopted

$$\pi_{\text{alg}}(a \mid s) = \frac{1}{2}\tilde{\beta}(a \mid s) + \frac{1}{2}\pi(a \mid s),$$

where π is the current policy and $\tilde{\beta}$ is defined as

$$\tilde{\beta}(\cdot \mid s) = \mathcal{N}(\boldsymbol{\mu}_{\hat{\beta}^*}(s), 4\boldsymbol{\Sigma}_{\hat{\beta}^*}(s)).$$

Here, $\mu_{\hat{\beta}^*}(s)$ and $\Sigma_{\hat{\beta}^*}(s)$ are the mean vector and covariance matrix of the selected behaviour policy $\hat{\beta}^*$. The first term regularises the critic values over a broad range of actions to guide a randomly initialised network $g_{\psi_{\pi}}$ towards producing near-zero values. The second term regularises the critic values in the vicinity of the current policy allowing delicate control near the boundary of OOD(s).

D DIDACTIC EXPERIMENT

To show the importance of trajectory clustering, we also evaluated our algorithm on a simple navigation environment shown in Figure 4. The agent should output a two-dimensional velocity vector

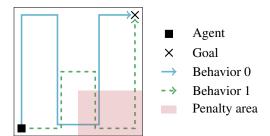


Figure 4: The agent receives a reward of -1 until it reaches the goal and an additional penalty of -1 when it is inside the penalty area. The offline dataset consists of trajectories sampled using two behaviour policies β_0 (Behaviour 0) and β_1 (Behaviour 1).

as an action based on the two-dimensional position vector given as state. The episode terminates when the agent reaches the goal, and until then, it receives a reward of -1. When the agent is inside the penalty area, it receives an additional penalty of -1, that is, the reward is -2 every time-step. We created an offline dataset that consists of trajectories sampled using two behaviour policies. Table 5 shows the mean and standard deviation of episode returns for three different algorithms: TD3+BC (Fujimoto & Gu, 2021), SVR (Mao et al., 2023), and ours. Due to the multi-modality of the dataset, SVR, which relies on simple behaviour cloning, fails to learn an optimal policy. Our algorithm can successfully recover the behaviour policies and thus outperforms other baselines.

D.1 IMPLEMENTATION DETAILS

We normalised the observations following Fujimoto & Gu (2021) and scaled the rewards following Kostrikov et al. (2022). Note that we used the naive behaviour cloning algorithm for expert and medium datasets, instead of our trajectory clustering method, because we know a priori that they are homogeneous. The algorithm was implemented upon the JAX (Bradbury et al., 2018) framework using the Flax (Heek et al., 2024) library. The scikit-learn (Pedregosa et al., 2011) library was used to compute ARIs and NMIs for Sections 5.3 and E.3.

E EXPERIMENT DETAILS

E.1 DIDACTIC EXPERIMENTS

The observation space is $S = [0, 30] \times [0, 30]$, the action space is $A = [-0.2, 0.2] \times [-0.2, 0.2]$, and the penalty area is $S_p = [15, 30] \times [0, 10]$. The starting location of the agent is sampled uniformly at random from $[0, 0.1] \times [0, 0.1]$ and the goal position is fixed to g = (30, 30). If the current state is $s = (s_0, s_1) \in S$ and the action is $a = (a_0, a_1) \in A$, the next state $s' = (s'_0, s'_1) \in S$ is

$$s' = \text{clip}(s + a; 0, 30).$$

The reward function is

$$r(s,a,s') = \begin{cases} 0 & \text{if } \|s'-g\|_2 < 0.1, \\ -2 & \text{if } s' \in \mathcal{S}_p, \\ -1 & \text{otherwise.} \end{cases}$$

Table 5: Average returns on the navigation dataset.

Algorithm	Return			
TD3+BC	-622.82 ± 185.72			
SVR	-1000.00 ± 0.00			
Ours	-173.50 ± 25.32			

1297

1298

1299

1300 1301

1302

1303

1304

1305 1306

1326 1327

1328

1330 1331

1332

1333

1334

1335 1336 1337

1338 1339

1340 1341

1342

1343

1344

1345

1347

1348

1349

The episode terminates if either the agent has approached the goal ($||s'-g||_2 < 0.1$) or the number of time-steps exceeded 1000.

The offline RL datasets was generated using subgoal-reaching policies. A subgoal-reaching policy $\pi(s; g_s)$ for a subgoal g_s is defined as

$$\pi(s; g_s) = \text{clip}(\text{clip}(g_s - s; -0.2, 0.2) + 0.1\varepsilon; -0.2, 0.2), \tag{15}$$

were ε is a two-dimensional standard Gaussian noise. We generated the samples according to Algorithm 1 using two different list of subgoals: [(0,30),(10,30),(10,0),(20,0),(20,30)] and [(10,0),(10,15),(20,15),(20,0),(30,0)].

Algorithm 1 Dataset generation from a list of subgoals

```
1307
             1: Input: a list of subgoals [g_s^{(1)}, g_s^{(2)}, \dots, g_s^{(N)}]
1309
             2: Initialize an empty dataset \mathcal{D}
             3: while \mathcal{D} has less than 1 000 000 elements do
1310
             4:
                       s \leftarrow \texttt{env.reset} ()
1311
                       for g_s \leftarrow [g_s^{(1)}, g_s^{(2)}, \dots, g_s^{(N)}] do while \|s - g_s\|_2 \ge 0.1 do
             5:
1312
             6:
1313
             7:
                                  a \leftarrow \pi(s; g_s)
                                                                                                                                                 ▷ (15)
1314
             8:
                                  s', r, d \leftarrow \text{env.step}(a)
1315
             9:
                                  Add (s, a, r, d) to \mathcal{D}
1316
                                  s \leftarrow s'
            10:
1317
                            end while
            11:
1318
            12:
                       end for
1319
            13:
                       while ||s - g||_2 \ge 0.1 do
1320
            14:
                             a \leftarrow \pi(s;g)
                                                                                                                                                 \triangleright (15)
1321
                             s', r, d \leftarrow \text{env.step}(a)
            15:
1322
            16:
                             Add (s, a, r, d) to \mathcal{D}
            17:
                             s \leftarrow s'
1323
                       end while
            18:
1324
            19: end while
1325
```

E.2 MOTIVATION EXPERIMENT

E.3 ADDITIONAL ANALYSIS ON THE TRAJECTORY CLUSTERING ALGORITHM

Aside from the three random-medium-expert datasets mentioned in Section 5.3, we also created custom D4RL datasets by concatenating medium and expert datasets of halfcheetah, hopper, and walker2d tasks. The mean and standard deviation of ARIs and NMIs for each configuration over 5 different seeds are reported in Table 6. The configuration $(\lambda_T, \lambda_R) = (1,0)$ performs the best on average even after we include the three medium-expert datasets.

F ADDITIONAL FIGURES

G RELATED WORK

Value regularisation Value regularisation aims to discourage the actor from choosing OOD actions by penalising their critic values. Conservative Q learning (CQL; Kumar et al. 2020) was one of the first works in this line of research, where they minimise the standard TD error together with the Q-values of OOD actions. Lyu et al. (2022) pointed out that the CQL excessively regularises the OOD Q-values to the extent that hampers the learning process. They suggested a milder regularisation term based on the critic values of ID actions. Supported value regularisation (SVR; Mao et al. 2023) proposed a penalisation scheme that maintains standard Bellman updates for ID actions while selectively penalising OOD actions' critic values. Most existing value regularisation algorithm, including the three works introduced in this section, sample the OOD actions from the current policy. However, as training progresses, the current policy will start to produce ID samples, so it is crucial to

Table 6: The impact of hyperparameters λ_T and λ_R on the performance of our trajectory clustering algorithm evaluated on custom D4RL datasets. The performance is measured in terms of adjusted rand index (ARI) and normalised mutual information score (NMI).

		λ_T =	= 1	$\lambda_T = 0$		
		$\lambda_R = 1$	$\lambda_R = 0$	$\lambda_R = 1$	$\lambda_R = 0$	
halfcheetah-med	lium-ex	pert				
	ARI	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	1.00 ± 0.00	
	NMI	1.00 ± 0.00	0.99 ± 0.01	1.00 ± 0.00	1.00 ± 0.00	
halfcheetah-rand	lom-me	dium-expert				
	ARI	0.97 ± 0.04	0.99 ± 0.00	0.91 ± 0.19	0.91 ± 0.17	
	NMI	0.97 ± 0.03	0.97 ± 0.01	0.93 ± 0.12	0.92 ± 0.12	
hopper-medium-	-expert					
	ARI	0.80 ± 0.44	1.00 ± 0.00	0.99 ± 0.01	1.00 ± 0.00	
	NMI	0.80 ± 0.42	1.00 ± 0.00	0.97 ± 0.03	1.00 ± 0.00	
hopper-random-	mediun	n-expert				
	ARI	0.49 ± 0.29	0.98 ± 0.02	0.59 ± 0.33	0.97 ± 0.06	
	NMI	0.57 ± 0.26	0.97 ± 0.02	0.66 ± 0.29	0.98 ± 0.04	
walker2d-mediu	m-expe	rt				
	ARI	0.99 ± 0.01	1.00 ± 0.00	0.99 ± 0.02	1.00 ± 0.00	
	NMI	0.99 ± 0.02	1.00 ± 0.01	0.98 ± 0.04	1.00 ± 0.00	
walker2d-randor	m-medi	um-expert				
	ARI	0.88 ± 0.16	0.98 ± 0.05	0.98 ± 0.03	1.00 ± 0.00	
	NMI	0.90 ± 0.11	0.98 ± 0.03	0.97 ± 0.04	1.00 ± 0.00	
Average	ARI	0.86 ± 0.27	0.99 ± 0.02	0.91 ± 0.21	0.98 ± 0.07	
	NMI	0.87 ± 0.24	0.98 ± 0.02	0.92 ± 0.17	0.98 ± 0.06	

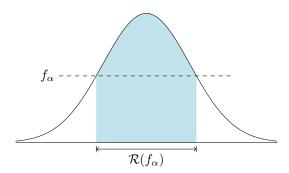


Figure 5: A diagram showing the probability density function and the $100(1-\alpha)$ % highest density region of a normal distribution. The probability of the corresponding normal random variable to lie inside $\mathcal{R}(f_{\alpha})$, which corresponds to the area of the coloured region, is $1-\alpha$.

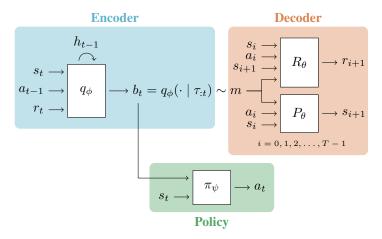


Figure 6: VariBAD architecture. This figure is a redrawn version of Figure 2 in Zintgraf et al. (2021).

prevent unnecessary penalisation for those actions. CQL circumvents this issue through maximising the critic values for actions in the dataset. SVR does it by soft thresholding the regulariser based on the importance sampling ratio. In contrary, our method adopts a hard thresholding mechanism where ID actions are not penalised at all. This is possible due to our capability of explicitly identifying the OOD action set.

Heterogeneous datasets There are multiple prior work concerned with offline RL datasets with heterogeneous behaviours. Wang et al. (2023) utilises a diffusion model (Sohl-Dickstein et al., 2015; Ho et al., 2020) to capture the multi-modality of the true behaviour policy. Li et al. (2023) trains a mixture of Gaussian policy on the dataset via likelihood maximisation and then obtains a closed-form estimate of the best possible action near the behaviour policy. These two works ignores the trajectory information and handles each transition individually. Mao et al. (2024) incorporates an expectation—maximisation algorithm to learn diverse policies from a given offline RL dataset. Wang et al. (2024) proposes a learning-based trajectory clustering algorithm that can also automatically determine the cluster size. Although these two works leverage the trajectory information, they obtain the trajectory representation by simply averaging the samples, causing a substantial loss of information. We incorporate a sequence modelling technique instead to learn an effective representation of each trajectory.

VQ-VAE State-conditioned action quantisation (SAQ; Luo et al. 2023) is closely related to our work in the sense that they also leverage a VQ-VAE in the offline RL setting. However, their main focus is to discretise the actions because most of the challenges in offline RL originates from the ambiguity of continuous distributions. On the other hand, our algorithm uses VQ-VAE to cluster

the trajectories and recover the behaviour policies. Also, SAQ discretises the actions individually, ignoring the trajectory information.

H LIMITATIONS

Our work is built upon the assumption that each trajectory in the dataset was sampled from a single behaviour policy. Although this assumption does not hold in general, as the behaviour policy may change mid-trajectory, the change is subtle enough for our algorithm to perform reasonably well. However, this might not be the case for real world scenarios. Future work could explore mechanisms to detect behaviour policy change and split the trajectory at those transition points.