

A Generalization of the definition

Adapted from Bäuerle and Rieder (2010), a more general definition of an episodic CMOMDP with horizon $H \in \mathbb{N}$ is given by a sequence of sets of data $\{(\mathcal{S}, \mathcal{A}, \mathcal{X}_h, \text{Pr}_h, c_h)\}_{h=1}^H$, where

- \mathcal{S} is the state space equipped with a σ -algebra $\mathcal{F}_{\mathcal{S}}$.
- \mathcal{A} is the action space equipped with a σ -algebra $\mathcal{F}_{\mathcal{A}}$.
- \mathcal{X}_h is a measurable subset of $\mathcal{X} = \mathcal{S} \times \mathcal{A}$, which denotes the set of accessible state-action pairs at time h .
- Pr_h is a stochastic transition kernel at time h . For any fixed $(s, a) \in \mathcal{X}_h$, the mapping $S \mapsto \text{Pr}_h(S | s, a)$ is a probability measure on \mathcal{S} . Moreover, the mapping $(s, a) \mapsto \text{Pr}_h(S | s, a)$ is measurable with respect to (s, a) for all $S \in \mathcal{F}_{\mathcal{S}}$. Intuitively, $\text{Pr}_h(S | s, a)$ gives the probability that the next state is in S if the current state is s and action a is taken at time h .
- $c_h : \mathcal{X}_h \rightarrow \mathbb{R}^D$ is a measurable function, which defines the expected cost at time h if the current state is s and action a is taken.

For any $h \in [H]$ and $(s, a) \in \mathcal{S} \times \mathcal{A}$, we define $\mathcal{P}_h(\cdot | s, a)$ as the density of $\text{Pr}_h(\cdot | s, a)$ relative to a certain measure μ by the Radon–Nikodym derivative, i.e., $\mathcal{P}_h = d\text{Pr}_h/d\mu$. We usually let μ be the Lebesgue measure when the state is continuous and be the counting measure when \mathcal{S} is discrete. For example, for measurable real-valued function f , when \mathcal{S} is discrete, $\mathcal{P}_h(\cdot | s, a)$ is the probability mass function and $\int_{\mathcal{S}} f d\mu$ coincides with $\sum_{s \in \mathcal{S}} f(s)$. When \mathcal{S} is continuous, $\mathcal{P}_h(\cdot | s, a)$ is the probability density function and $\int_{\mathcal{S}} f d\mu$ denotes the Lebesgue integration. Our proposed method is able to handle both of the above situations. For notational simplicity, we omit the dependence on μ through the paper.

B Explanations and justifications of Assumption 1

In this section, we explain and justify Assumption 1 imposed on the target set. This assumption is a stronger version of the one used in Yu et al. (2021) where \mathcal{W}^* is not assumed to be a lower set.

B.1 Nonsingular intersections

We first explain the meaning of nonsingular intersection, namely the existence of γ_{\max} . The intuition is that the target set \mathcal{W}^* and the achievable values \mathcal{V} do not share the same support hyperplane.

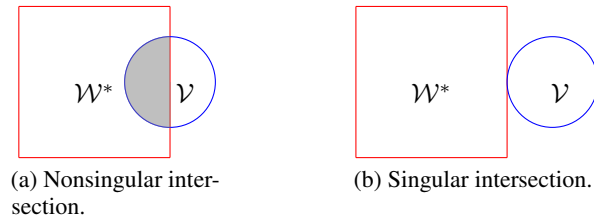


Figure 1: The illustration of two types of intersections: the nonsingular intersection and the singular intersection. The target set \mathcal{W}^* is represented by a red square and the achievable values \mathcal{V} is represented by a blue circle. The intersection is in grey.

We illustrate it via a two-dimensional example in Figure 1. The target set \mathcal{W}^* is represented by a red square and the achievable values \mathcal{V} is represented by a blue circle. The intersection is in grey. In (a), the intersection is nonsingular, which guarantees the existence of $\gamma_{\max} < \pi$. However, in (b), \mathcal{V} intersects with \mathcal{W} singularly, leading to the inexistence of an upper bound γ_{\max} . We notice that Figure 1 also intuitively shows why the Slater’s condition implies our assumption, since interior points guarantee the nonsingular intersection.

As we see in related proofs, for any $\mathbf{W} \in \mathcal{V}$, nonsingular intersection enables us to reduce the calculation of $\text{dist}(\mathbf{W}, \mathcal{W}^*)$ to $\text{dist}(\mathbf{W}, \mathcal{W})$ multiplied by the sine of the included angle. We note

that when this geometric requirement is not satisfied (i.e., when the intersection is singular), this reduction may not exist. The reason is that when \mathcal{V} is not of full dimension, \mathbf{W} can be arbitrarily close to \mathcal{W}^* while remains far away from \mathcal{W} , making the reduction impossible (Yu et al., 2021).

B.2 Being a lower set

Now we discuss the assumption of being a lower set. It is much milder than at first glance and can be easily relaxed. For target set \mathcal{W}^* which is not a lower set, we can transform it into a lower set with at most $D \cdot 2^D$ dimensions. We achieve this by cutting the boundary of \mathcal{W}^* into at most 2^D pieces. We provide a two-dimensional example as shown in Figure 2. The target set is defined as

$$\mathcal{W}^* = \{(c^1, c^2) \mid c^1 + c^2 \geq H/2, c^1 + c^2 \leq 3H/2, c^1 - c^2 \geq -H/2, c^1 - c^2 \leq H/2\}.$$

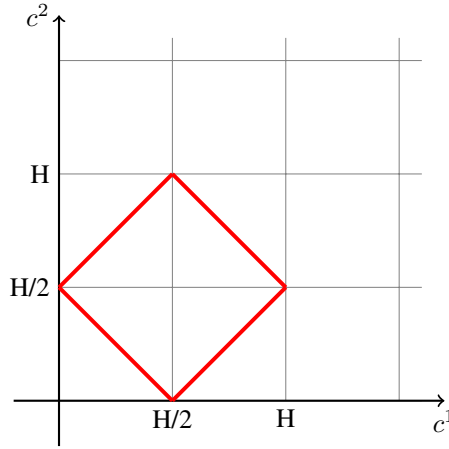


Figure 2: An example of target set $\mathcal{W}^* = \{(c^1, c^2) \mid c^1 + c^2 \geq H/2, c^1 + c^2 \leq 3H/2, c^1 - c^2 \geq -H/2, c^1 - c^2 \leq H/2\}$. Although it does not satisfy Assumption 1 of being a lower set, we can cut its boundaries into four pieces, each of which corresponds to a partial constraint. Since each partial constraint can be expressed in a way that is a lower set, the target set \mathcal{W}^* can be reformulated as the combination of those partial constraints to satisfy the lower set requirement at the expense of higher dimension.

We cut its boundaries into four pieces: the upper left one, the upper right one, the bottom left one, and the bottom right one. For the upper right boundary, it corresponds to the partial constraint:

$$\{(c^1, c^2) \mid c^1 + c^2 \leq 3H/2\}.$$

It is a little tricky to represent the other three. We take the upper left boundary, $c^2 - c^1 \leq H/2$, for an example. To that end, we notice that it is equivalent to $c^2 + (H - c^1) \leq 3H/2$. Hence, we set $c^3 = H - c^1$ and $c^4 = c^2$. Then we can represent this boundary by $c^3 + c^4 \leq 3H/2$, which satisfies Assumption 1 of being a lower set with respect to c^3 and c^4 . Analogously, we introduce c_5 to c_8 for the other two boundaries. Finally, we obtain a new target set with 8 objectives, c_1, \dots, c_8 , that meet Assumption 1. Hence, we successfully relax the constraint of being a lower set.

B.3 Summaries

In a word, our assumption is mild, and thus target sets that satisfy Assumption 1 are very common in applications. For example, we can define \mathcal{W}^* as $\mathcal{W}^* = \{\mathbf{W} \in [0, H]^D : \mathbf{W} \leq \mathbf{b}\}$ for $\mathbf{b} \in \mathbb{R}^D$, which is equivalent to the formation of constraints in CMDPs.

C Reducing CMOMDPs to CMDPs

C.1 Proof of Proposition 1

Proof of Proposition 1. We first show that the CMDP is a special case of the CMOMDP. Then, we point out the correspondence between the assumptions imposed on CMDPs and on CMOMDPs. To be specific, we show that Slater’s condition imposed on CMDPs is equivalent to the existence of γ_{\max} in Assumption 1. Hence, the CMOMDP is a reasonable generalization of CMDPs.

Reducing to CMDPs. An episodic CMDP is usually given by $\mathcal{M} = (\mathcal{S}, \mathcal{A}, H, \mathcal{P}, c, \mathbf{u})$, where \mathcal{S} is a state space, \mathcal{A} is an action space, H is the horizon, $\mathcal{P} = \{\mathcal{P}_h\}_{h=1}^H$ is a collection of transition kernels $\mathcal{P}_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, $c = \{c_h\}_{h=1}^H$ is a collection of cost functions $c_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, and $\mathbf{u} = \{\mathbf{u}_h\}_{h=1}^H$ is a collection of utility functions $\mathbf{u}_h : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]^D$.

The CMDP aims to solve the following problem:

$$\min_{\pi \in \Delta(\mathcal{A} | \mathcal{S}, H)} V_1^{c, \pi}(\underline{s}) \quad \text{s.t.} \quad \mathbf{V}_1^{\mathbf{u}, \pi}(\underline{s}) \leq \mathbf{b}, \quad (11)$$

where $V^{c, \pi}$ and $\mathbf{V}^{\mathbf{u}, \pi}$ are the state-value functions for cost c and utility \mathbf{u} , respectively. Here we assume $\mathbf{b} \in [0, H]^D$ to avoid triviality.

Now we construct a CMOMDP \mathcal{M}' that is equivalent to \mathcal{M} . To that end, we set $\mathcal{M}' = (\mathcal{S}, \mathcal{A}, H, \mathcal{P}, \mathbf{c}')$ where $\mathcal{S}, \mathcal{A}, H$, and \mathcal{P} are same as \mathcal{M} . We define $\mathbf{c}' = \{c'_h\}_{h=1}^H$ as the concatenation of c and \mathbf{u} , i.e., $c'_h = (c_h, \mathbf{u}_h^\top)^\top \in [0, 1]^{D+1}$ for $h \in [H]$.

The preference function $g : \mathbb{R}^{D+1} \rightarrow \mathbb{R}$ is defined as the first coordinate map, i.e., we define $g(x_1, \dots, x_{D+1}) = x_1$.

We set the target set \mathcal{W}^* as $\mathcal{W}^* = \{\mathbf{W} \in [0, H]^{D+1} : \mathbf{W}_{2:(D+1)} \leq \mathbf{b}\}$ where $\mathbf{W}_{2:(D+1)}$ denotes the D -dimensional vector obtained by removing the first coordinate of \mathbf{W} .

The CMOMDP \mathcal{M}' seeks to solve the following optimization problem

$$\min_{\pi \in \Delta(\mathcal{A} | \mathcal{S}, H)} g(V_1^{\mathbf{c}', \pi}(\underline{s})) \quad \text{s.t.} \quad \mathbf{V}_1^{\mathbf{c}', \pi}(\underline{s}) \in \mathcal{W}^*. \quad (12)$$

A careful analysis will reveal that, by the construction of \mathcal{M}' , (12) is equivalent to (11). Hence, CMDP is a special case of CMOMDP.

Correspondence of assumptions. Under the above definition, it is clear that g is 1-Lipschitz and convex and that $g(\mathbf{x}) \geq g(\mathbf{x}')$ as long as $\mathbf{x} \geq \mathbf{x}'$ for any $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{D+1}$. In addition, \mathcal{W}^* is a lower set, and it is close and convex.

It remains to show that the existence of γ_{\max} in Assumption 1 is satisfied. However, this is equivalent to Slater’s condition on the CMDPs. Recall that Slater’s condition supposes the existence of an interior point, which is equivalent to a nonsingular intersection between the target set and the set of achievable values. Hence, it is exactly the geometric analog of Slater’s condition. \square

C.2 PEDI on CMDPs

Now we analyze PEDI on CMDPs. Recall that the core of PEDI is solving (7). Based on the analysis in Appendix C.1, as we have specified g as the first coordinate map, we can derive its convex conjugate,

$$g^*(\mathbf{x}^*) = \sup_{\mathbf{x}} (\mathbf{x}^\top \mathbf{x}^* - g(\mathbf{x})) = \sup_{\mathbf{x}} (\mathbf{x}^\top \mathbf{x}^* - x_1) = \begin{cases} 0, & \mathbf{x}^* = (1, 0, \dots, 0)^\top \\ +\infty, & \text{otherwise} \end{cases}. \quad (13)$$

Therefore, solving (7) is equivalent to solving

$$\mathbf{p}^* = \min_{V \in \mathcal{V}} \max_{\alpha \in \mathcal{B}^D} V^1 - V^{1,*} + \nu \alpha^\top \mathbf{V} - \nu \max_{\mathbf{x} \in \mathcal{W}^*} \alpha^\top \mathbf{x}, \quad (14)$$

since the objective attains the maximum with respect to β when $\beta = (1, 0, \dots, 0)^\top$. Here, by V^1 and $V^{1,*}$ we mean the first coordinates of \mathbf{V} and \mathbf{V}^* , respectively. Now we consider $\max_{\mathbf{x} \in \mathcal{W}^*} \alpha^\top \mathbf{x}$. Recall that we construct \mathcal{W}^* by

$$\mathcal{W}^* = \{\mathbf{W} \in [0, H]^{D+1} : \mathbf{W} \leq \mathbf{b}'\}, \quad (15)$$

where we suppose $\mathbf{b}' = (H, \mathbf{b}^\top)^\top$ for simplicity. We observe this constraint is element-wise, i.e., different coordinates of \mathbf{W} are independent. Hence, we have

$$\max_{\mathbf{x} \in \mathcal{W}^*} \boldsymbol{\alpha}^\top \mathbf{x} = \max_{\mathbf{x} \in \mathcal{W}^*} \sum_{i=1}^{D+1} \alpha_i x_i = \max_{\substack{0 \leq x_i \leq b_i, \\ i \in [D+1]}} \sum_{i=1}^{D+1} \alpha_i x_i = \sum_{i=1}^{D+1} \max_{\substack{0 \leq x_i \leq b_i, \\ i \in [D+1]}} \alpha_i x_i = \boldsymbol{\alpha}_{i,+}^\top \mathbf{b}', \quad (16)$$

where $\boldsymbol{\alpha}_{i,+}^\top$ denotes $(\boldsymbol{\alpha}_i^\top)_+$. Then, we have

$$\mathbf{p}^* = \min_{\mathbf{V} \in \mathcal{V}} \max_{\boldsymbol{\alpha} \in \mathcal{B}^D} V^1 - V^{1,*} + \nu(\boldsymbol{\alpha}^\top \mathbf{V} - \boldsymbol{\alpha}_+^\top \mathbf{b}'). \quad (17)$$

A good observation is that $\boldsymbol{\alpha} \geq 0$ always holds (otherwise, we can increase the negative component of $\boldsymbol{\alpha}$ to make the objective larger). Therefore, the optimization problem is actually

$$\begin{aligned} \mathbf{p}^* &= \min_{\mathbf{V} \in \mathcal{V}} \max_{\boldsymbol{\alpha} \in \mathcal{B}_+^D} V^1 - V^{1,*} + \nu \boldsymbol{\alpha}(\mathbf{V} - \mathbf{b}') \\ &= \min_{\mathbf{V} \in \mathcal{V}} \max_{\boldsymbol{\alpha} \in \mathcal{B}_+^D} V_1^{c,\pi}(\underline{s}) - V_1^{c,*}(\underline{s}) + \nu \boldsymbol{\alpha}(\mathbf{V}_1^{u,\pi}(\underline{s}) - \mathbf{b}), \end{aligned} \quad (18)$$

where the second equality is by definition and the fact that $\mathbf{b}'_1 = H \geq V^1$. We notice that $\nu \boldsymbol{\alpha}$ can be considered the Lagrangian multiplier. Hence, we conclude that the objective reduces to the Lagrangian formulation of the CMDP optimization problem, and our algorithm reduces to the dual gradient method.

D Equivalence of problems

In this section, we show that (1) and (4) share the same solution.

For convenience, we denote by \mathbf{p}_1^* the original problem (1), and by \mathbf{p}_2^* the new problem (4). We first substitute their variables π with the state-value function.

For \mathbf{p}_1^* , we have

$$\begin{aligned} \mathbf{p}_1^* &= \min_{\pi} g(\mathbf{V}_1^\pi(\underline{s})) \quad \text{s.t.} \quad \mathbf{V}_1^\pi(\underline{s}) \in \mathcal{W}^* \\ &= \min_{\mathbf{V} \in \mathcal{V}} g(\mathbf{V}) \quad \text{s.t.} \quad \mathbf{V} \in \mathcal{W}^* \end{aligned} \quad (19)$$

while for \mathbf{p}_2^* , we have

$$\begin{aligned} \mathbf{p}_2^* &= \min_{\pi} (\text{SubOpt}(\pi) + \nu \text{Violation}(\pi)) \\ &= \min_{\pi} (g(\mathbf{V}_1^\pi(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \nu \text{dist}(\mathbf{V}_1^\pi(\underline{s}), \mathcal{W}^*)) \\ &= \min_{\mathbf{V} \in \mathcal{V}} (g(\mathbf{V}) - g(\mathbf{V}_1^*(\underline{s})) + \nu \text{dist}(\mathbf{V}, \mathcal{W}^*)) \\ &= \min_{\mathbf{V} \in \mathcal{V}} g(\mathbf{V}) + \nu \text{dist}(\mathbf{V}, \mathcal{W}^*) \end{aligned} \quad (20)$$

We suppose that \mathbf{V}^\dagger is the solution of \mathbf{p}_2^* . To show that \mathbf{V}^\dagger is also the solution of \mathbf{p}_1^* , it suffices to verify $\text{dist}(\mathbf{V}^\dagger, \mathcal{W}^*) = 0$.

Suppose $\text{dist}(\mathbf{V}^\dagger, \mathcal{W}^*) \neq 0$, i.e., $\mathbf{V}^\dagger \notin \mathcal{W}^*$. We consider its projection on \mathcal{W}^* , $\mathbf{V}^\S = \prod_{\mathcal{W}^*} \mathbf{V}^\dagger$. Since g is 1-Lipschitz, when $\nu > 1$, it holds that

$$g(\mathbf{V}^\S) - g(\mathbf{V}^\dagger) < \nu \text{dist}(\mathbf{V}^\dagger, \mathbf{V}^\S). \quad (21)$$

Then, by (21), we have

$$\begin{aligned} &g(\mathbf{V}^\S) + \nu \text{dist}(\mathbf{V}^\S, \mathcal{W}^*) \\ &< g(\mathbf{V}^\S) - g(\mathbf{V}^\dagger) + g(\mathbf{V}^\dagger) + \nu(\text{dist}(\mathbf{V}^\S, \mathcal{W}^*) + \text{dist}(\mathbf{V}^\dagger, \mathbf{V}^\S)) \\ &= g(\mathbf{V}^\dagger) + \nu \text{dist}(\mathbf{V}^\dagger, \mathcal{W}^*), \end{aligned}$$

where the last equality is for $\text{dist}(\mathbf{V}^\S, \mathcal{W}^*) = 0$. It contradicts our assumption that \mathbf{V}^\dagger is the solution of \mathbf{p}_2^* . Therefore, we conclude that the solution \mathbf{V}^\dagger of \mathbf{p}_2^* should lie in \mathcal{W}^* , which completes the proof of the equivalence between (1) and (4).

E An instantiation of PEDI for linear kernel CMOMDPs

This section serves as complements for Section 3.3.

We now propose a method to estimate the empirical transition kernel $\hat{\mathcal{P}}$ and cost function \hat{c} and thereby construct a ξ -uncertainty quantifier to specify Algorithm 2 for linear kernel CMOMDPs. We construct the empirical transition kernel by ridge regression on the offline dataset \mathcal{D} as follows,

$$\begin{aligned} \hat{\mathcal{P}}_h(s' | s, a) &= \psi(s, a, s')^\top \hat{\theta}_h \\ \text{where } \hat{\theta}_h &= \arg \min_{\theta \in \mathbb{R}^{d_1}} \sum_{\tau=1}^N \int_{\mathcal{S}} (\delta_{s_{h+1}^\tau}(s') - \psi(s_h^\tau, a_h^\tau, s')^\top \theta)^2 ds' + \lambda \|\theta\|_2^2. \end{aligned}$$

Here $\delta_{s_h^\tau}(s)$ is the Dirac function centered at s_h^τ for continuous space and indicator function for finite space, and $\lambda > 0$ is the regularization parameter. Note that we can obtain the following closed form of $\hat{\theta}_h$,

$$\hat{\theta}_h = \Lambda_h^{-1} \sum_{\tau=1}^N \psi(s_h^\tau, a_h^\tau, s_{h+1}^\tau) \quad \text{where} \quad \Lambda_h = \sum_{\tau=1}^N \int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \psi(s_h^\tau, a_h^\tau, s')^\top ds' + \lambda I. \quad (22)$$

We construct \hat{c} in an analogous way by ridge regression,

$$\hat{r}_h^i(s, a) = \varphi(s, a)^\top \hat{\theta}_h^{c^i} \quad \text{where} \quad \hat{\theta}_h^{c^i} = \arg \min_{\theta \in \mathbb{R}^{d_2}} \sum_{\tau=1}^N (c_h^{i,\tau} - \varphi(s_h^\tau, a_h^\tau)^\top \theta)^2 + \lambda \|\theta\|_2^2.$$

and $\hat{\theta}_h^{c^i}$ has the closed form

$$\hat{\theta}_h^{c^i} = \Lambda_{\varphi,h}^{-1} \sum_{\tau=1}^N c_h^{i,\tau} \cdot \varphi(s_h^\tau, a_h^\tau) \quad \text{where} \quad \Lambda_{\varphi,h} = \sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau) \varphi(s_h^\tau, a_h^\tau)^\top + \lambda I. \quad (23)$$

Moreover, we construct the ξ -uncertainty quantifier for $h \in [H]$ below

$$\Gamma_h^{\mathcal{P}}(s, a, s') = \min \{ \kappa \cdot \|\psi(s, a, s')\|_{\Lambda_h^{-1}}, 1 \}, \quad \Gamma_h^{c^i}(s, a) = \min \{ \kappa \cdot \|\varphi(s, a)\|_{\Lambda_{\varphi,h}^{-1}}, 1 \}, \quad (24)$$

where $\kappa > 0$ is a scaling parameter to be specified later. By plugging (24) into the pessimistic planning (Algorithm 1), we finish the establishment of PEDI on linear kernel CMOMDPs.

F Pessimism is all you need

This section studies the effectiveness of the pessimistic approach for offline CMOMDPs. To that end, we first introduce the model evaluation error and then develop the decomposition lemma, which decomposes the discrepancy between the state-value functions of the learned policy and the optimal policy into three parts: the spurious correlation, the intrinsic uncertainty, and the optimization error. Then, we show that our proposed method successfully eliminates the spurious correlation, which is the most difficult one to control.

We consider a meta-algorithm that constructs a sequence of policies $\{\pi_k\}_{k=1}^K$ which ideally converges to the optimal policy. At the k -th iteration, the algorithm constructs the estimations $\{V_h^k\}_{h=1}^H$ and $\{Q_h^k\}_{h=1}^H$ such that $V_h^k(s) = \mathbb{D}_{\pi^k}[Q_h^k](s)$. We define the model evaluation error below, which characterizes the error of estimating the Bellman equations.

Definition 3 (Model evaluation error). *The model evaluation error for the k -th iteration is defined as*

$$\iota_h^k(s, a) = Q_h^k(s, a) - \mathbb{B}_h V_{h+1}^k(s, a).$$

We denote its i -th scalar component by $\iota_h^{i,k}(s, a)$, i.e., $\iota_h^{i,k}(s, a) = Q_h^{i,k}(s, a) - \mathbb{B}_h V_{h+1}^{i,k}(s, a)$.

Utilizing the model evaluation error, the discrepancy between any state-value function and the optimal one admits a decomposition as shown in the following lemma.

Lemma 2 (Decomposition of state-value function). *Let $\{\mathbf{V}_h^k\}_{h=1}^H$ and $\{\mathbf{Q}_h^k\}_{h=1}^H$ be any state-value function and action-value function such that $\mathbf{V}_h^k(s) = \mathbb{D}_{\pi^k}[\mathbf{Q}_h^k](s)$ for any $s \in \mathcal{S}$ and any $h \in [H]$. Then, it holds that*

$$\begin{aligned} \mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^*(\underline{s}) &= \sum_{h=1}^H \mathbb{E}_{\pi^*}[\boldsymbol{\iota}_h^k(s_h, a_h) \mid s_1 = \underline{s}] \\ &\quad + \sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\left\langle \mathbf{Q}_h^k(s_h, \cdot), (\pi_h^k - \pi_h^*)(\cdot \mid s_h) \right\rangle_{\mathcal{A}} \mid s_1 = \underline{s} \right] \end{aligned}$$

and that

$$\begin{aligned} \mathbf{V}_1^{\pi^k}(\underline{s}) - \mathbf{V}_1^*(\underline{s}) &= - \underbrace{\sum_{h=1}^H \mathbb{E}_{\pi^k}[\boldsymbol{\iota}_h^k(s_h, a_h) \mid s_1 = \underline{s}]}_{\text{(i) Spurious Correlation}} + \underbrace{\sum_{h=1}^H \mathbb{E}_{\pi^*}[\boldsymbol{\iota}_h^k(s_h, a_h) \mid s_1 = \underline{s}]}_{\text{(ii) Intrinsic Uncertainty}} \\ &\quad + \underbrace{\sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\left\langle \mathbf{Q}_h^k(s_h, \cdot), (\pi_h^k - \pi_h^*)(\cdot \mid s_h) \right\rangle_{\mathcal{A}} \mid s_1 = \underline{s} \right]}_{\text{(iii) Optimization Error}}, \end{aligned}$$

where \mathbb{E}_{π^k} and \mathbb{E}_{π^*} are taken with respect to the trajectories induced by π^k and π^* in the underlying CMOMDP, respectively.

Proof of Lemma 2. See Appendix H.6.1 for a detailed proof. \square

Lemma 2 is the vectorized analogue of Lemma 3.1 in Jin et al. (2020b). It suggests that we can decompose the discrepancy between the state-value function of learned policy and the optimal one into (i) spurious correlation, (ii) intrinsic uncertainty, and (iii) optimization error. Among them, (i) is the most difficult to control since it depends on both π^k and $\boldsymbol{\iota}_h^k$ that spuriously correlated with each other. As the learner has no control over the data collecting process, this spurious correlation could be large even in a multi-armed bandit setting (Jin et al., 2020b). Term (ii) is easier to control since π^* is intrinsic to the underlying CMOMDP and therefore not spuriously correlated with $\boldsymbol{\iota}_h^k$.

As proved in Section 4.1, our proposed algorithm successfully eliminates term (i) through pessimism. In Section 4.2, we have shown that (ii) is impossible to eliminate as it arises from the information-theoretic lower bound of linear kernel CMOMDPs.

G Reducing linear kernel CMOMDPs to tabular CMOMDPs

All we need is to represent the transition kernel \mathcal{P} and cost function c of tabular CMOMDP in the form of linear kernel CMOMDP.

We set $d_1 = |\mathcal{S}||\mathcal{A}||\mathcal{S}|$ and $d_2 = |\mathcal{S}||\mathcal{A}|$ and set $\psi(s, a, s') = \mathbf{e}_{(s,a,s')}, (\theta_h)_{(s,a,s')} = \mathcal{P}(s' \mid s, a), \varphi(s, a) = \mathbf{e}_{(s,a)}$, and $(\theta_h^{c_i})_{(s,a)} = c_h^i(s, a)$.

Here we denote by \mathbf{e} the canonical basis. It can be verified that the definition of linear kernel MDP (Definition 2) is satisfied with $R = 1$.

H Proofs for Section 4

For notational simplicity, we sometime use the shorthand $\boldsymbol{\iota}_h^k, \mathbf{Q}_h, \mathbf{V}_h, \pi_h$ to denote $\boldsymbol{\iota}_h^k(s_h, a_h), \mathbf{Q}_h(s_h, a_h), \mathbf{V}_h(s_h)$, and $\pi_h(a_h \mid s_h)$ when there is no risk of confusion.

H.1 Proof of Lemma 1

Proof of Lemma 1. The proof is by backward induction. Suppose the inequality holds for Q -values in the $(h+1)$ -th step. Then, it holds that for any $s \in \mathcal{S}, i \in [D]$,

$$V_{h+1}^{i,k}(s) = \mathbb{D}_{\pi_h^k}[\mathbf{Q}_{h+1}^{i,k}](s) \geq \mathbb{D}_{\pi_h^k}[\mathbf{Q}_{h+1}^{i,\pi^k}](s) = V_{h+1}^{i,\pi^k}.$$

For any $(s, a) \in \mathcal{S} \times \mathcal{A}$, $i \in [D]$, when $Q_h^{i,k}(s, a) = H - h + 1$, it holds that $Q_h^{i,k}(s, a) \geq Q_h^{i,\pi^k}(s, a)$. Otherwise, by the definition of ξ -uncertainty quantifier defined in Definition 1, it holds that

$$\begin{aligned} Q_h^{i,k}(s, a) &\geq \tilde{c}_h^i(s, a) + \hat{\mathbb{P}}_h^k[V_{h+1}^{i,k}](s, a) + \Gamma_h + \Gamma_h^{c^i} \\ &\geq c_h^i(s, a) + \mathbb{P}_h^k[V_{h+1}^{i,k}](s, a) \\ &\geq c_h^i(s, a) + \mathbb{P}_h^k[V_{h+1}^{i,\pi^k}](s, a) \\ &= Q_h^{i,\pi^k}(s, a) \end{aligned}$$

under event \mathcal{E} . Thus, by induction, it holds for all h that $Q_h^{i,k} \geq Q_h^{i,\pi^k}$.

Then, by the fact that $g(\mathbf{x}) \geq g(\mathbf{x}')$ holds as long as $\mathbf{x} \geq \mathbf{x}'$, we have

$$g(\hat{\mathbf{V}}_1(\underline{s})) \geq g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})).$$

For the constraint violation, we consider the point $\mathbf{W} = \mathbf{V}_1^{\pi^k}(\underline{s}) + \prod_{\mathcal{W}^*} \mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^k(\underline{s})$. By the fact that \mathcal{W}^* is a lower set in Assumption 1, it holds that $\mathbf{W}_+ \in \mathcal{W}^*$ and

$$\text{dist}(\mathbf{V}_1^{\pi^k}(\underline{s}), \mathcal{W}^*) \leq \text{dist}(\mathbf{V}_1^{\pi^k}(\underline{s}), \mathbf{W}) \leq \text{dist}(\mathbf{V}_1^k(\underline{s}), \mathcal{W}^*).$$

Thus, we complete the proof of Lemma 1. □

H.2 Proof of Theorem 1

Proof of Theorem 1. In the following lemma, we show that the difference between any state-value function and the optimal one, as stated in Lemma 2, can be bounded from above by the ξ -uncertainty quantifier with high probability when projected along with θ^k .

Lemma 3 (Upper bound of projected difference of state-value functions). *Suppose $\{(\Gamma_h^{\mathcal{P}}, \Gamma_h^{\mathcal{C}})\}_{h=1}^H$ in Algorithm 2 is a ξ -uncertainty quantifier. Then under event \mathcal{E} , we have*

$$(\theta^k)^\top (\mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^*(\underline{s})) \leq 2(1 + \rho)\sqrt{D} \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) + \|\Gamma_h^{\mathcal{C}}(s_h, a_h)\|_\infty \mid s_1 = \underline{s}]$$

Proof of Lemma 3. See Appendix H.6.3 for a detailed proof. □

In what follows, we suppose the event \mathcal{E} holds, which has probability at least $1 - \xi$. Applying Lemma 1, we have

$$\begin{aligned} &K \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \right] \\ &\leq K \left[g(\hat{\mathbf{V}}_1(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\hat{\mathbf{V}}_1(\underline{s}), \mathcal{W}^*) \right]. \end{aligned}$$

By the convex conjugate in (6) and the fact that $\hat{\mathbf{V}}_1(\underline{s}) = \sum_{k=1}^K \mathbf{V}_1^k(\underline{s})$, we get

$$\begin{aligned} &K \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \right] \\ &= \max_{\|\beta\| \leq 1} \left\{ \beta \cdot \sum_{k=1}^K \mathbf{V}_1^k(\underline{s}) - \sum_{k=1}^K g^*(\beta) \right\} - Kg(\mathbf{V}_1^*(\underline{s})) \\ &\quad + \rho \max_{\|\alpha\| \leq 1} \left\{ \alpha \cdot \sum_{k=1}^K \mathbf{V}_1^k(\underline{s}) - \sum_{k=1}^K \max_{\mathbf{x} \in \mathcal{W}^*} \alpha \cdot \mathbf{x} \right\} \end{aligned}$$

We observe these two terms

$$\max_{\|\beta\| \leq 1} \left\{ \beta \cdot \sum_{k=1}^K \mathbf{V}_1^k(\underline{s}) - \sum_{k=1}^K g^*(\beta) \right\}, \quad \max_{\|\alpha\| \leq 1} \left\{ \alpha \cdot \sum_{k=1}^K \mathbf{V}_1^k(\underline{s}) - \sum_{k=1}^K \max_{\mathbf{x} \in \mathcal{W}^*} \alpha \cdot \mathbf{x} \right\},$$

are the “single best desicion” in hindsight in the projected subgradient method. By setting $\eta^k = 2G^{-1}\sqrt{D/k}$ (or $2G^{-1}\sqrt{D/K}$ if K is predefined), we apply Theorem 5 with $R = 2$ and $G = 2(1 + \rho)H\sqrt{D}$ (to verify the conditions, note that g^* is $H\sqrt{D}$ -Lipschitz) to get

$$\begin{aligned} & K \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \right] \\ & \leq \sum_{k=1}^K \left\{ (\beta^k)^\top \mathbf{V}_1^k(\underline{s}) - g^*(\beta^k) \right\} - Kg(\mathbf{V}_1^*(\underline{s})) + \rho \sum_{k=1}^K \left\{ (\alpha^k)^\top \mathbf{V}_1^k(\underline{s}) - \max_{\mathbf{x} \in \mathcal{W}^*} (\alpha^k)^\top \mathbf{x} \right\} \\ & \quad + C(1 + \rho)\sqrt{DH^2K} \end{aligned} \quad (25)$$

where C is a constant. Then, by observing

$$g^*(\beta^k) = \max_{\mathbf{V}} \{ (\beta^k)^\top \mathbf{V} - g(\mathbf{V}) \} \geq (\beta^k)^\top \mathbf{V}_1^*(\underline{s}) - g(\mathbf{V}_1^*(\underline{s})),$$

$$\max_{\mathbf{x} \in \mathcal{W}^*} (\alpha^k)^\top \mathbf{x} \geq (\alpha^k)^\top \mathbf{V}_1^*(\underline{s}),$$

we have for (25) that

$$\begin{aligned} & K \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \right] \\ & \leq \sum_{k=1}^K \left\{ (\beta^k)^\top \mathbf{V}_1^k(\underline{s}) - (\beta^k)^\top \mathbf{V}_1^*(\underline{s}) + \rho((\alpha^k)^\top \mathbf{V}_1^k(\underline{s}) - (\alpha^k)^\top \mathbf{V}_1^*(\underline{s})) \right\} + C(1 + \rho)\sqrt{DH^2K} \\ & = \sum_{k=1}^K \left[(\theta^k)^\top (\mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^*(\underline{s})) \right] + C(1 + \rho)\sqrt{DH^2K} \end{aligned} \quad (26)$$

Applying Lemma 3, we have

$$\begin{aligned} & K \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \right] \\ & \leq 2K(1 + \rho)\sqrt{D} \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) + \|\Gamma_h^c(s_h, a_h)\|_\infty \mid s_1 = \underline{s}] + C(1 + \rho)\sqrt{DH^2K} \\ & = K(\epsilon_K + \text{IntUncert}_{\mathcal{D}}^{\pi^*}), \end{aligned}$$

where we define

$$\epsilon_K = C(1 + \rho)\sqrt{\frac{DH^2}{K}}, \quad \text{IntUncert}_{\mathcal{D}}^{\pi^*} = 2(1 + \rho)\sqrt{D} \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) + \|\Gamma_h^c(s_h, a_h)\|_\infty \mid s_1 = \underline{s}].$$

Note that $\text{dist}(\mathbf{W}^K, \mathcal{W}^*) \geq 0$. Therefore, we can bound the suboptimality from above by

$$g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) \leq \epsilon_K + \text{IntUncert}_{\mathcal{D}}^{\pi^*}.$$

To obtain an upper bound of the constraint violation, we employ the following lemma.

Lemma 4. *Let \mathbf{W}^* denote a return vector in set \mathcal{W} that achieves the lowest cost, i.e. $\forall \mathbf{W} \in \mathcal{W}, g(\mathbf{W}) \geq g(\mathbf{W}^*)$. Then, under Assumption 1, it holds for any $\mathbf{W} \in \mathbb{R}^D$ that*

$$g(\mathbf{W}) - g(\mathbf{W}^*) \geq -\text{dist}(\mathbf{W}, \mathcal{W}^*) / \sin(\gamma_{\max}).$$

Proof of Lemma 4. See Lemma 16 in Yu et al. (2021) for a detailed proof. \square

We notice that by definition, $\mathbf{W}^* = \mathbf{V}_1^*(\underline{s})$. Therefore, applying Lemma 4 with $\mathbf{W} = \mathbf{V}_1^{\hat{\pi}}(\underline{s})$, we have

$$\begin{aligned}
& \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \\
& \leq \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) + \sin(\gamma_{\max}) \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \frac{\text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*)}{\sin(\gamma_{\max})} \right] \\
& = \sin(\gamma_{\max}) \left[g(\mathbf{V}_1^{\hat{\pi}}(\underline{s})) - g(\mathbf{V}_1^*(\underline{s})) + \rho \text{dist}(\mathbf{V}_1^{\hat{\pi}}(\underline{s}), \mathcal{W}^*) \right] \\
& \leq \frac{2}{\rho} (\epsilon_K + \text{IntUncert}_{\mathcal{D}}^{\pi^*}).
\end{aligned}$$

where the last inequality follows from (26). Thus, we complete the proof of Theorem 1. \square

H.3 Proof of Theorem 2

Proof of Theorem 2. It suffices to show that $\{(\Gamma_h^{\mathcal{P}}, \Gamma_h^{\mathcal{C}})\}_{h=1}^H$ defined in (24) is a ξ -uncertainty quantifier as stated in the following lemma.

Lemma 5. *Under Assumptions 2, we set*

$$\lambda = 1, \quad \kappa = CR\sqrt{d \log(dN) + \log(DH/\xi)},$$

where $C > 0$ is an absolute constant and $\xi \in (0, 1)$ is the confidence parameter. Then, $\{(\Gamma_h^{\mathcal{P}}, \Gamma_h^{\mathcal{C}})\}_{h=1}^H$ in (24) is a ξ -uncertainty quantifier.

Proof of Lemma 5. See Appendix H.6.5 for a detailed proof. \square

By Theorem 1, we have that

$$\text{SubOpt}(\hat{\pi}) \leq \epsilon_K + \text{IntUncert}_{\mathcal{D}}^{\pi^*}, \quad \text{Violation}(\hat{\pi}) \leq \frac{2}{\rho} (\epsilon_K + \text{IntUncert}_{\mathcal{D}}^{\pi^*})$$

where we define

$$\epsilon_K = C(1 + \rho)\sqrt{\frac{DH^2}{K}}, \quad \text{IntUncert}_{\mathcal{D}}^{\pi^*} = 2(1 + \rho)\sqrt{D} \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s, a) + \|\Gamma_h^{\mathcal{C}}(s, a)\|_{\infty} \mid s_1 = \underline{s}].$$

By Lemma 5 and the ξ -uncertainty quantifier defined in (24), we finish the proof. \square

H.4 Proof of Theorem 3

Proof of Theorem 3. The following lemma is adopted from Jin et al. (2020b), which characterizes the information-theoretic lower bound of offline RL.

Theorem 4. *For the output $\hat{\pi}$ of any offline RL algorithm, there exists a tabular MDP \mathcal{M} with initial state $\underline{s} \in \mathcal{S}$ and a dataset \mathcal{D} compliant with \mathcal{M} , such that*

$$\mathbb{E}_{\mathcal{D}} \left[\frac{\text{SubOpt}(\hat{\pi})}{\sum_{h=1}^H \mathbb{E}_{\pi^*} [1/\sqrt{1 + n_h(s_h, a_h)} \mid s_1 = \underline{s}]} \right] \geq C$$

where $C > 0$ is an absolute constant. Here $n_h(s_h, a_h) = \sum_{\tau=1}^N \mathbb{1}\{s_h^{\tau} = s_h, a_h^{\tau} = a_h\}$ for $(s_h, a_h) \in \mathcal{S} \times \mathcal{A}$.

Proof of Theorem 4. See Theorem 4.6 in Jin et al. (2020b) for a detailed proof. \square

Note that we can view the MDP as a special case of CMOMDP with the target set $\mathcal{W}^* = \mathbb{R}^D$, $D = 1$ and g being the identity function. Hence the hard instance in Theorem 4 is also a hard instance of CMOMDP. It remains to reduce the tabular MDP there to a linear kernel MDP defined in Definition 2. To that end, we set $d_1 = |\mathcal{S}||\mathcal{A}||\mathcal{S}|$ and $d_2 = |\mathcal{S}||\mathcal{A}|$ and set $\psi(s, a, s') = \mathbf{e}_{(s,a,s')}, (\theta_h)_{(s,a,s')} =$

$\mathcal{P}(s' | s, a), \varphi(s, a) = \mathbf{e}_{(s,a)}$, and $(\theta_h^{c_i})_{(s,a)} = c_h^i(s, a)$. Here \mathbf{e} denotes the canonical basis. It can be verified that Definition 2 is satisfied with $R = 1$.

Then, it holds that

$$\Lambda_h = \lambda I + \sum_{\tau=1}^N \sum_{s' \in \mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \psi(s_h^\tau, a_h^\tau, s')^\top = \lambda I + \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \sum_{s' \in \mathcal{S}} n_h(s, a) E_{(s,a,s'), (s,a,s')}$$

where $E_{(s,a,s'), (s,a,s')}$ is the matrix in which entries at $((s, a, s'), (s, a, s'))$ is 1 and other entries are all 0. We note that $\lambda = 1$ and Λ_h is diagonal and thus we have

$$\|\psi(s, a, s')\|_{\Lambda_h^{-1}} \leq \frac{1}{\sqrt{1 + n_h(s, a)}} \quad (27)$$

Following the same derivation we get

$$\|\varphi(s, a)\|_{\Lambda_{\varphi,h}^{-1}} \leq \frac{1}{\sqrt{1 + n_h(s, a)}} \quad (28)$$

Then, by Theorem 4, (27), and (28), we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} \left[\frac{\text{SubOpt}(\hat{\pi})}{\sum_{h=1}^H \mathbb{E}_{\pi^*} [\|\varphi(s_h, a_h)\|_{\Lambda_{\varphi,h}^{-1}} + |\mathcal{S}|^{-1} \int_{\mathcal{S}} \|\psi(s_h, a_h, s')\|_{\Lambda_h^{-1}} ds' \mid s_1 = \underline{s}]} \right] \\ & \geq \mathbb{E}_{\mathcal{D}} \left[\frac{\text{SubOpt}(\hat{\pi})}{2 \sum_{h=1}^H \mathbb{E}_{\pi^*} [1/\sqrt{1 + n_h(s_h, a_h)} \mid s_1 = \underline{s}]} \right] \\ & \geq c/2. \end{aligned}$$

which completes the proof of Theorem 3. \square

H.5 Proof of Corollary 1

Proof of Corollary 1. By the property of visitation measure, we have that

$$\begin{aligned} & \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) + \Gamma_h^{c_i}(s_h, a_h) \mid s_1 = \underline{s}] \\ & = \mathbb{E}_{\mu_h^*} [\Gamma_h(s, a)] + \mathbb{E}_{\mu_h^*} [\Gamma_h^{c_i}(s, a)] \\ & \leq \left(\sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} \frac{\mu_h^*(s, a)}{\mu_{h,b,\tau}(s, a)} \right) \cdot \left(\mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h(s, a)] + \mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h^{c_i}(s, a)] \right) \\ & \leq c \cdot \left(\mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h(s, a)] + \mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h^{c_i}(s, a)] \right), \end{aligned} \quad (29)$$

where the first inequality follows from Hölder's inequality and the last inequality follows from the condition. We will upper bound the term $\mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h(s, a)]$ and $\mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h^{c_i}(s, a)]$ respectively.

Let $X_h^{c_i, \tau} = \mathbb{E}_{\mu_{h,b,\tau}} [\Gamma_h^{c_i}(s, a)] - \Gamma_h^{c_i}(s_h^\tau, a_h^\tau)$, which is a martingale difference process with respect to the filtration $\{\mathcal{F}_h^\tau\}_{\tau=1}^N$. To see this, we have

$$\mathbb{E}[X_h^{c_i, \tau} \mid \mathcal{F}_h^{\tau-1}] = \mathbb{E}[X_h^{c_i, \tau}] = 0$$

Note that $|X_h^{c_i, \tau}|$ is bounded by 2, and thus Azuma's inequality implies for all $h \in [H]$ that

$$\sum_{\tau=1}^N X_h^{c_i, \tau} \leq C \sqrt{N \log(DH/\xi)} \quad (30)$$

holds with probability at least $1 - \xi/(D + 1)$. Here C denotes an absolute constant.

Moreover, by the Cauchy-Schwarz inequality, it holds that

$$\sum_{\tau=1}^N \Gamma_h^{c_i}(s_h^\tau, a_h^\tau) \leq \sum_{\tau=1}^N \kappa \sqrt{\varphi(s_h^\tau, a_h^\tau)^\top \Lambda_{\varphi,h}^{-1} \varphi(s_h^\tau, a_h^\tau)} \leq \kappa \sqrt{N \sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau)^\top \Lambda_{\varphi,h}^{-1} \varphi(s_h^\tau, a_h^\tau)}, \quad (31)$$

and by the property of trace, we have

$$\begin{aligned}
& \sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau)^\top \Lambda_{\varphi,h}^{-1} \varphi(s_h^\tau, a_h^\tau) \\
&= \text{tr} \left(\sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau)^\top \Lambda_{\varphi,h}^{-1} \varphi(s_h^\tau, a_h^\tau) \right) \\
&= \text{tr} \left(\sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau) \varphi(s_h^\tau, a_h^\tau)^\top \Lambda_{\varphi,h}^{-1} \right) \\
&= \text{tr} \left((\Lambda_{\varphi,h} - \lambda I) \Lambda_{\varphi,h}^{-1} \right) \\
&= \text{tr} \left((I - \lambda \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_{d_2})^{-1}) \right) \\
&\leq d_2,
\end{aligned} \tag{32}$$

where in the last equality we denote by $\lambda_1, \lambda_2, \dots, \lambda_{d_2}$ the eigenvalues of $\Lambda_{\varphi,h}$. By plugging (32) into (31) we have

$$\sum_{\tau=1}^N \Gamma_h^{c_i}(s_h^\tau, a_h^\tau) \leq \kappa \sqrt{dN}. \tag{33}$$

Combining (30) and (33), we have

$$\sum_{\tau=1}^N \mathbb{E}_{\mu_h^{\text{b},\tau}} [\Gamma_h^{c_i}(s, a)] = \sum_{\tau=1}^N X_h^{c_i,\tau} + \sum_{\tau=1}^N \Gamma_h^{c_i}(s_h^\tau, a_h^\tau) \leq C\kappa \sqrt{dN \log(DH/\xi)} \tag{34}$$

with probability at least $1 - \xi/(D+1)$.

For term $\mathbb{E}_{\mu_h^{\text{b},\tau}} [\Gamma_h(s, a)]$, we follow a similar derivation. Let $X_h^\tau = \mathbb{E}_{\mu_h^{\text{b},\tau}} [\Gamma_h(s, a)] - \Gamma_h(s_h^\tau, a_h^\tau)$ be a martingale difference sequence with bound $2H$ for each X_h^τ . Then, by the Azuma's inequality, it holds that

$$\sum_{\tau=1}^N X_h^\tau \leq CH \sqrt{N \log(DH/\xi)} \tag{35}$$

with probability at least $1 - \xi/(D+1)$. Similarly, we have

$$\begin{aligned}
\sum_{\tau=1}^N \Gamma_h(s_h^\tau, a_h^\tau) &\leq \sum_{\tau=1}^N H \int_{\mathcal{S}} \kappa \sqrt{\psi(s_h^\tau, a_h^\tau, s')^\top \Lambda_h^{-1} \psi(s_h^\tau, a_h^\tau, s')} \, ds' \\
&\leq \kappa H \sqrt{|\mathcal{S}|N \sum_{\tau=1}^N \int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s')^\top \Lambda_h^{-1} \psi(s_h^\tau, a_h^\tau, s') \, ds'}
\end{aligned} \tag{36}$$

and note that

$$\begin{aligned}
& \sum_{\tau=1}^N \int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s')^\top \Lambda_h^{-1} \psi(s_h^\tau, a_h^\tau, s') \, ds' \\
&= \text{tr} \left(\sum_{\tau=1}^N \int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s')^\top \Lambda_h^{-1} \psi(s_h^\tau, a_h^\tau, s') \, ds' \right) \\
&= \text{tr} \left(\sum_{\tau=1}^N \int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \psi(s_h^\tau, a_h^\tau, s')^\top \Lambda_h^{-1} \, ds' \right) \\
&\leq d_1.
\end{aligned} \tag{37}$$

Combining (35), (36) and (37), we have

$$\sum_{\tau=1}^N \mathbb{E}_{\mu_h^{\text{b},\tau}} [\Gamma_h(s, a)] = \sum_{\tau=1}^N X_h^\tau + \sum_{\tau=1}^N \Gamma_h(s_h^\tau, a_h^\tau) \leq C\kappa H \sqrt{dN|\mathcal{S}| \log(DH/\xi)} \tag{38}$$

holds with probability at least $1 - \xi/(D + 1)$.

Taking the union bound for (34) and (38), by (29), we have

$$\begin{aligned} \text{IntUncert } \bar{\pi}_D^* &= 2(1 + \rho) \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) + \|\Gamma_h^c(s_h, a_h)\|_\infty \mid s_1 = \underline{s}] \\ &\leq 2\varsigma(1 + \rho) \sum_{h=1}^H N^{-1} \sum_{\tau=1}^N \left(\mathbb{E}_{\mu_h^{b,\tau}} [\Gamma_h(s, a)] + \mathbb{E}_{\mu_h^{b,\tau}} [\|\Gamma_h^c(s, a)\|_\infty] \right) \\ &\leq C\varsigma\kappa(1 + \rho)H^2\sqrt{d|\mathcal{S}|/N\log(DH/\xi)} \end{aligned}$$

with probability at least $1 - \xi$. Thus, we complete the proof of Corollary 1. \square

H.6 Supporting lemmas and proofs

H.6.1 Proof of Lemma 2

Proof of Lemma 2.

Lemma 6 (Extended value difference (Cai et al., 2020)). *Let $\pi = \{\pi_h\}_{h=1}^H$ and $\pi' = \{\pi'_h\}_{h=1}^H$ be two arbitrary policies and let $\{Q_h\}_{h=1}^H$ be any given Q -functions. For any $h \in [H]$, we define a value function $V_h : \mathcal{S} \rightarrow \mathbb{R}$ by letting $V_h(s) = \langle Q_h(s, \cdot), \pi_h(\cdot | s) \rangle_{\mathcal{A}}$ for all $s \in \mathcal{S}$. Then, we have*

$$\begin{aligned} V_1(\underline{s}) - V_1^{\pi'}(\underline{s}) &= \sum_{h=1}^H \mathbb{E}_{\pi'} [\langle Q_h(s_h, \cdot), \pi_h(\cdot | s_h) - \pi'_h(\cdot | s_h) \rangle_{\mathcal{A}} \mid s_1 = \underline{s}] \\ &\quad + \sum_{h=1}^H \mathbb{E}_{\pi'} [Q_h(s_h, a_h) - (\mathbb{B}_h V_{h+1})(s_h, a_h) \mid s_1 = \underline{s}]. \end{aligned}$$

where \underline{s} is an initial state.

Proof of Lemma 6. See Section B.1 in Cai et al. (2020) for a detailed proof. \square

Applying Lemma 6 with $\pi = \pi^k$, $\pi' = \pi^*$, and $Q_h = Q_h^{i,k}$, we get

$$V_1^{i,k}(\underline{s}) - V_1^{i,*}(\underline{s}) = \sum_{h=1}^H \mathbb{E}_{\pi^*} [\langle Q_h^{i,k}, \pi_h^k - \pi_h^* \rangle \mid s_1 = \underline{s}] + \sum_{h=1}^H \mathbb{E}_{\pi^*} [\ell_h^{i,k} \mid s_1 = \underline{s}]. \quad (39)$$

Moreover, applying Lemma 6 with $\pi = \pi' = \pi^k$, we get

$$V_1^{i,\pi^k}(\underline{s}) - V_1^{i,k}(\underline{s}) = - \sum_{h=1}^H \mathbb{E}_{\pi^k} [\ell_h^{i,k} \mid s_1 = \underline{s}]. \quad (40)$$

Since for any $k \in [K]$,

$$\mathbf{V}_1^{\pi^k}(\underline{s}) - \mathbf{V}_1^*(\underline{s}) = \mathbf{V}_1^{\pi^k}(\underline{s}) - \mathbf{V}_1^k(\underline{s}) + \mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^*(\underline{s}), \quad (41)$$

we are done by plugging (39) and (40) into (41). \square

H.6.2 Projected subgradient method

Our algorithm benefits from the online projected subgradient method for the update of dual variables. We formally state it below for compactness.

Online learning. Online learning involves two players: the adversary and the player. The online learning protocol is shown in Algorithm 3.

Algorithm 3 Protocol of Online Learning

```

1: for  $t = 1, \dots, T$  do
2:   The player chooses an action  $x_t$ .
3:   The adversary picks a function  $f_t$ .
4:   The player obtains reward  $f_t(x_t)$ .
5:   The player learns via  $f_t$ .
6: end for

```

Note that there is no assumption on how the adversary will pick the function f_t , and it may be adversarially chosen. The player aims to minimize the regret:

$$\text{Regret} = \max_x \sum_{t=1}^T f_t(x) - \sum_{t=1}^T f_t(x_t), \quad (42)$$

which measures the quality of the player's strategy x_1, \dots, x_T compared with the single best decision in hindsight.

Projected subgradient method. The projected subgradient method is a particular case of mirror descent/ascent with Euclidean distance. Applying this method to online learning produces a regret bound of the order $O(\sqrt{T})$.

We suppose that the actions x_t are required to be contained in some convex set \mathcal{X} , i.e., $x_t \in \mathcal{X}$. Let $g_t \in \partial f_t(x_t)$ denote a subgradient of f_t at x_t and G and R denote two constant bounds such that $\max_{x, y \in \mathcal{X}} \|x - y\|_2 \leq R$ and $\max_{t \in [T]} \|\partial f_t(x_t)\|_2 \leq G$. We set the step length η_t at the t -th iteration to $\frac{R}{G\sqrt{t}}$ if we do not know the number of iterations T in advance and to $\frac{R}{G\sqrt{T}}$ if we have the knowledge of T . The latter case will lead to an upper bound with a smaller constant multiplicative factor. With these notations, the update rule of projected subgradient method can be expressed as

$$x_{t+1} \leftarrow \arg \max_{x \in \mathcal{X}} \left\{ f_t(x_t) + \langle \eta_t g_t, x - x_t \rangle - \frac{1}{2} \|x - x_t\|_2^2 \right\}.$$

We describe the complete method in Algorithm 4.

Algorithm 4 projected subgradient method

```

1: Arbitrarily initialize  $x_1 \in \mathcal{X}$ .
2: for  $t = 1, \dots, T - 1$  do
3:    $x_{t+1} \leftarrow \arg \max_{x \in \mathcal{X}} \left\{ f_t(x_t) + \langle \eta_t g_t, x - x_t \rangle - \frac{1}{2} \|x - x_t\|_2^2 \right\}$ 
4: end for

```

By this method, the regret is guaranteed to increase sublinearly as stated in Theorem 5.

Theorem 5. *Using projected subgradient method mentioned in Algorithm 4, it holds that for the regret (42) that*

$$\text{Regret} \leq CRG\sqrt{T},$$

where C is an absolute constant.

Proof of Theorem 5. See Zinkevich (2003) for a detailed proof. □

H.6.3 Proof of Lemma 3

Proof of Lemma 3. The model evaluation error can be upper bounded by the ξ -uncertainty quantifier. We formally state it below.

Lemma 7. *Under event \mathcal{E} , it holds that for any $k \in [K]$ and $h \in [H]$,*

$$0 \leq \iota_h^k(s, a) \leq 2(\Gamma_h(s, a) \cdot \mathbf{1} + \Gamma_h^c(s, a)).$$

Proof of Lemma 7. See Appendix H.6.4 for a detailed proof. □

By Lemma 2, it holds that

$$\begin{aligned}
(\boldsymbol{\theta}^k)^\top (\mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^*(\underline{s})) &= \underbrace{\sum_{h=1}^H \mathbb{E}_{\pi^*} \left[(\boldsymbol{\theta}^k)^\top \boldsymbol{\iota}_h^k \mid s_1 = \underline{s} \right]}_{\text{(i)}} \\
&\quad + \underbrace{\sum_{h=1}^H \mathbb{E}_{\pi^*} \left[\langle (\boldsymbol{\theta}^k)^\top \mathbf{Q}_h^k, \pi_h^k - \pi_h^* \rangle_{\mathcal{A}} \mid s_1 = \underline{s} \right]}_{\text{(ii)}}.
\end{aligned}$$

We bound these two terms above separately. For (ii), since π^k in Algorithm 2 is greedy (see Line 5), we have (ii) ≤ 0 . For (i), by Lemma 7, we get

$$(i) \leq 2|\boldsymbol{\theta}^k|^\top \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) \cdot \mathbf{1} + \Gamma_h^c(s_h, a_h) \mid s_1 = \underline{s}].$$

By plugging them back and applying Hölder's inequality, we have

$$(\boldsymbol{\theta}^k)^\top (\mathbf{V}_1^k(\underline{s}) - \mathbf{V}_1^*(\underline{s})) \leq 2(1 + \rho)\sqrt{D} \sum_{h=1}^H \mathbb{E}_{\pi^*} [\Gamma_h(s_h, a_h) + \|\Gamma_h^c(s_h, a_h)\|_\infty \mid s_1 = \underline{s}],$$

where we notice that $\|\boldsymbol{\theta}^k\|_1 \leq (1 + \rho)\sqrt{D}$. Thus, we finish the proof of Lemma 3. \square

H.6.4 Proof of Lemma 7

Proof of Lemma 7. For any $i \in [D]$, recall that we have

$$Q_h^{i,k}(s, a) = \min \{ \overline{Q}_h^{i,k}(s, a), H - h + 1 \}_+.$$

Under event \mathcal{E} , we have

$$\begin{aligned}
\overline{Q}_h^{i,k}(s, a) &= \widehat{c}_h^i(s, a) + \widehat{\mathbb{P}}_h[V_{h+1}^{i,k}](s, a) + \Gamma_h + \Gamma_h^c \\
&\geq c_h^i(s, a) + \mathbb{P}_h[V_{h+1}^{i,k}] \\
&\geq 0,
\end{aligned}$$

where the last inequality follows from $V_{h+1}^{i,k} \in [0, H - h]$. Therefore, it holds that $Q_h^{i,k}(s, a) \leq \overline{Q}_h^{i,k}(s, a)$ and

$$\begin{aligned}
Q_h^{i,k}(s, a) &= \min \{ \overline{Q}_h^{i,k}(s, a), H - h + 1 \}_+ \\
&\geq \min \{ c_h^i(s, a) + \mathbb{P}_h[V_{h+1}^{i,k}], H - h + 1 \}_+ \\
&= c_h^i(s, a) + \mathbb{P}_h[V_{h+1}^{i,k}],
\end{aligned}$$

which implies

$$\iota_h^{i,k}(s, a) = Q_h^{i,k}(s, a) - [c_h^i(s, a) + \mathbb{P}_h V_{h+1}^{i,k}(s, a)] \geq 0.$$

It remains to establish an upper bound for $\iota_h^{i,k}(s, a)$. To that end, we have

$$\begin{aligned}
\iota_h^{i,k}(s, a) &= Q_h^{i,k}(s, a) - [c_h^i(s, a) + \mathbb{P}_h V_{h+1}^{i,k}(s, a)] \\
&\leq \overline{Q}_h^{i,k}(s, a) - [c_h^i(s, a) + \mathbb{P}_h V_{h+1}^{i,k}(s, a)] \\
&= [\widehat{c}_h^i(s, a) + \widehat{\mathbb{P}}_h[V_{h+1}^{i,k}](s, a) + \Gamma_h(s, a) + \Gamma_h^c(s, a)] - [c_h^i(s, a) + \mathbb{P}_h[V_{h+1}^{i,k}](s, a)] \\
&= [\Gamma_h^c(s, a) - c_h^i(s, a) + \widehat{c}_h^i(s, a)] + [\Gamma_h(s, a) - \mathbb{P}_h V_{h+1}^{i,k}(s, a) + \widehat{\mathbb{P}}_h[V_{h+1}^{i,k}](s, a)] \\
&\leq 2(\Gamma_h(s, a) + \Gamma_h^c(s, a)),
\end{aligned}$$

where the last inequality follows from the definition of the ξ -uncertainty quantifier. Thus, we finish the proof of Lemma 7. \square

H.6.5 Proof of Lemma 5

Proof of Lemma 5. In what follows, we show that $\{(\Gamma_h^{\mathcal{P}}, \Gamma_h^{\mathcal{C}})\}_{h=1}^H$ defined in (24) are ξ -uncertainty quantifier for linear kernel CMOMDP.

Uncertainty quantifier for \mathcal{P} . We first show that $\Gamma_h^{\mathcal{P}}$ is the ξ -uncertainty quantifier of \mathcal{P} .

By definition, we have

$$\begin{aligned}\mathcal{P}_h(s' | s, a) &= \psi(s, a, s')^\top \theta_h = \psi(s, a, s')^\top \Lambda_h^{-1} \Lambda_h \theta_h \\ &= \psi(s, a, s')^\top \Lambda_h^{-1} \left(\sum_{\tau=1}^N \int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \mathcal{P}_h(s' | s_h^\tau, a_h^\tau) ds' + \lambda \theta_h \right).\end{aligned}$$

Thus, by the closed form of $\hat{\theta}_h$ in (22), we have

$$\begin{aligned}\mathcal{P}_h(s' | s, a) - \hat{\mathcal{P}}_h(s' | s, a) &= \mathcal{P}_h(s' | s, a) - \psi(s, a, s')^\top \hat{\theta}_h \\ &= \underbrace{\psi(s, a, s')^\top \Lambda_h^{-1} \left(\sum_{\tau=1}^N \left(\int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \mathcal{P}_h(s' | s_h^\tau, a_h^\tau) ds' - \psi(s_h^\tau, a_h^\tau, s_{h+1}^\tau) \right) \right)}_{(i)} \\ &\quad + \underbrace{\lambda \cdot \psi(s, a, s')^\top \Lambda_h^{-1} \theta_h}_{(ii)}.\end{aligned}\tag{43}$$

For term (i), by Cauchy-Schwartz inequality, we have

$$\begin{aligned} |(i)| &\leq \|\psi(s, a, s')\|_{\Lambda_h^{-1}} \cdot \left\| \sum_{\tau=1}^N \left(\int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \mathcal{P}_h(s' | s_h^\tau, a_h^\tau) ds' - \psi(s_h^\tau, a_h^\tau, s_{h+1}^\tau) \right) \right\|_{\Lambda_h^{-1}} \\ &\leq C_1 R \cdot \sqrt{d \log(dN) + \log(DH/\xi)} \cdot \|\psi(s, a, s')\|_{\Lambda_h^{-1}}, \quad \forall (s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}, h \in [H]\end{aligned}$$

with probability at least $1 - \xi/(D+1)$. Now we prove the last inequality of the above derivation. To that end, we need the following lemma that generalizes the Theorem 1 in Abbasi-Yadkori et al. (2011) to function-valued process.

Lemma 8 (Self-normalized bound for function-valued Process). *Let Ω be a probability space and $\{\eta_t\}_{t=1}^\infty$ be a function-valued stochastic process with a filtration $\{\mathcal{G}_t\}_{t=0}^\infty$, i.e., $\eta_t : \Omega \times \mathcal{S} \rightarrow \mathbb{R}$. We assume that $\eta_t | \mathcal{G}_{t-1}$ is zero-mean and σ -sub-Gaussian, i.e.,*

$$\mathbb{E}[\eta_t(s) | \mathcal{G}_{t-1}] = 0, \quad \forall s \in \mathcal{S},$$

$$\log \mathbb{E} \left[\exp(\langle f, \eta_t \rangle) | \mathcal{G}_{t-1} \right] \leq \|f\|_\infty^2 \cdot \sigma^2 / 2, \quad \forall f : \mathcal{S} \rightarrow \mathbb{R}.$$

Let $\{X_t\}_{t=0}^\infty$ be an \mathbb{R}^d -function-valued stochastic process, i.e., $X_t : \Omega \times \mathcal{S} \rightarrow \mathbb{R}^d$, and suppose X_t is \mathcal{G}_{t-1} -measurable. We further assume that

$$\|\lambda^\top X_t\|_\infty \leq R \cdot \|\lambda^\top X_t\|_2 \tag{44}$$

almost surely for any $\lambda \in \mathbb{R}^d$. Let $V \in \mathbb{R}^{d \times d}$ be a positive definite matrix. We define

$$\bar{V}_t = V + \sum_{\tau=1}^t \int_{\mathcal{S}} X_\tau(s) X_\tau(s)^\top ds$$

and

$$S_t = \sum_{\tau=1}^t \langle X_\tau, \eta_\tau \rangle_{\mathcal{S}}.$$

Then for any $\delta > 0$, with probability at least $1 - \delta$, it holds that for any $t \geq 0$,

$$\|S_t\|_{\bar{V}_t^{-1}}^2 \leq 2\sigma^2 R^2 \log \left(\frac{\det(\bar{V}_t)^{1/2}}{\delta \det(V)^{1/2}} \right).$$

Proof of Lemma 8. See Appendix H.6.6 for a detailed proof. \square

We consider the filtration $\{\mathcal{F}_h^\tau\}_{h,\tau=1}^{H,N}$ defined in Assumption 2. Then it holds that

$$\mathbb{E}[\delta_{s_{h+1}^\tau}(s') | \mathcal{F}_h^\tau] = \mathcal{P}_h(s' | s_h^\tau, a_h^\tau).$$

For any $f : \mathcal{S} \rightarrow \mathbb{R}$, by Hölder's inequality, it holds that $\langle f, \mathcal{P}_h(\cdot | s_h^\tau, a_h^\tau) - \delta_{s_{h+1}^\tau} \rangle_{\mathcal{S}} \leq 2\|f\|_\infty$, which implies

$$\log \mathbb{E} \left[\exp \left(\langle f, \mathcal{P}_h(\cdot | s_h^\tau, a_h^\tau) - \delta_{s_{h+1}^\tau} \rangle \right) \middle| \mathcal{F}_{h,\tau} \right] \leq 4\|f\|_\infty^2/2.$$

It corresponds to the conditional 2-sub-Gaussianity. Moreover, noticing that $\psi(s_h^\tau, a_h^\tau, s')$ is \mathcal{F}_h^τ -measurable and both $\mathcal{P}(\cdot | s_h^\tau, a_h^\tau)$ and $\delta_{s_{h+1}^\tau}$ are \mathcal{F}_{h+1}^τ -measurable, we apply Lemma 8 with $X_\tau = \psi(s_h^\tau, a_h^\tau, \cdot)$, $\eta_h = \mathcal{P}_h(\cdot | s_h^\tau, a_h^\tau) - \delta_{s_{h+1}^\tau}$ and $V = \lambda I$ to get

$$\begin{aligned} & \left\| \sum_{\tau=1}^N \left(\int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \mathcal{P}_h(s' | s_h^\tau, a_h^\tau) ds' - \psi(s_h^\tau, a_h^\tau, s_{h+1}^\tau) \right) \right\|_{\Lambda_h^{-1}}^2 \\ & \leq 8R^2 \cdot \log(H/p \cdot \det(\Lambda_h)^{1/2} \det(\lambda I)^{-1/2}) \end{aligned} \quad (45)$$

with probability at least $1 - p/H$. It remains to upper bound $\det(\Lambda_h)$.

By Definition 2, we have

$$y^\top \Lambda_h y = \lambda \|y\|_2^2 + \sum_{\tau=1}^N \langle y^\top \psi(s_h^\tau, a_h^\tau, \cdot), y^\top \psi(s_h^\tau, a_h^\tau, \cdot) \rangle \leq \lambda \cdot \|y\|_2^2 + dN \cdot \|y\|_2^2,$$

which implies $\|\Lambda_h\|_2 \leq \lambda + dN$, and therefore,

$$\det(\Lambda_h) \leq \|\Lambda_h\|_2^d \leq (\lambda + dN)^d. \quad (46)$$

Setting $\lambda = 1$ and plugging (46) back into (45), we get

$$\begin{aligned} & \left\| \sum_{\tau=1}^N \left(\int_{\mathcal{S}} \psi(s_h^\tau, a_h^\tau, s') \mathcal{P}_h(s' | s_h^\tau, a_h^\tau) ds' - \psi(s_h^\tau, a_h^\tau, s_{h+1}^\tau) \right) \right\|_{\Lambda_h^{-1}}^2 \\ & \leq 8R^2 \cdot (1/2 \cdot d \log(1 + dN) + \log(H/p)) \\ & \leq CR^2 \cdot (d \log(dN) + \log(H/p)) \end{aligned} \quad (47)$$

holds with probability at least $1 - p/H$. Here C is an absolute constant. By the union bound, (47) holds for all $h \in [H]$ with probability at least $1 - p$.

For term (ii) in (43), by setting $\lambda = 1$, we have

$$|(ii)| \leq \|\psi(s, a, s')\|_{\Lambda_h^{-1}} \cdot \|\theta_h\|_{\Lambda_h^{-1}} \leq \sqrt{d} \cdot \|\psi(s, a, s')\|_{\Lambda_h^{-1}} \quad (48)$$

where the last inequality is due to the definition of linear kernel MDP (Definition 2) and $\|\Lambda_h^{-1}\|_2 \leq 1$.

By plugging (47) and (48) into (43), we get, for all $h \in [H]$, $(s, a, s') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$,

$$\begin{aligned} |\mathcal{P}_h(s' | s, a) - \widehat{\mathcal{P}}_h(s' | s, a)| & \leq CR \sqrt{d \log(dN) + \log(DH/\xi)} \cdot \|\psi(s, a, s')\|_{\Lambda_h^{-1}} \\ & \leq \kappa \cdot \|\psi(s, a, s')\|_{\Lambda_h^{-1}} \end{aligned} \quad (49)$$

holds with probability at least $1 - \xi/(D + 1)$.

Uncertainty quantifier for c . Due to the closed form solution of $\hat{\theta}_h^{c^i}$ in (23), we have

$$\begin{aligned}
& |c_h^i(s, a) - \hat{c}_h^i(s, a)| \\
&= |\varphi(s, a)^\top (\theta_h^{c^i} - \hat{\theta}_h^{c^i})| \\
&= |\varphi(s, a)^\top \Lambda_{\varphi, h}^{-1} \left(\sum_{\tau=1}^N c_h^{i, \tau} \cdot \varphi(s_h^\tau, a_h^\tau) - \Lambda_{\varphi, h} \theta_h^{c^i} \right)| \\
&= |\varphi(s, a)^\top \Lambda_{\varphi, h}^{-1} \left(\sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau) (c_h^{i, \tau} - \varphi(s_h^\tau, a_h^\tau)^\top \theta_h^{c^i}) - \lambda \theta_h^{c^i} \right)| \\
&\leq \|\varphi(s, a)\|_{\Lambda_{\varphi, h}^{-1}} \cdot \left\| \sum_{\tau=1}^N \varphi(s_h^\tau, a_h^\tau) (c_h^{i, \tau} - \varphi(s_h^\tau, a_h^\tau)^\top \theta_h^{c^i}) \right\|_{\Lambda_{\varphi, h}^{-1}} + \lambda \|\theta_h^{c^i}\|_{\Lambda_{\varphi, h}^{-1}}
\end{aligned}$$

Following a similar argument as we did for \mathcal{P} , we obtain a result analogous to (49),

$$\|c_h^i(s, a) - \hat{c}_h^i(s, a)\| \leq \kappa \cdot \|\varphi(s, a)\|_{\Lambda_{\phi, h}^{-1}}.$$

with probability at least $1 - p$.

Finally, by setting $p = \xi/(D + 1)$ and taking union bound for \mathcal{P} and c^i ($i \in [D]$), we complete the proof of Lemma 5. \square

H.6.6 Proof of Lemma 8

Proof of Lemma 8. We generalize the proof of Abbasi-Yadkori et al. (2011) as follows.

Lemma 9. *Let $\lambda \in \mathbb{R}^d$ be an arbitrary vector and*

$$M_t^\lambda = \exp \left[\sum_{\tau=1}^t \left(\frac{\langle \lambda^\top X_\tau, \eta_\tau \rangle}{\sigma^2 R^2} - \frac{\|\lambda^\top X_\tau\|_2^2}{2} \right) \right].$$

Let T be a stopping time with respect to $\{\mathcal{G}_t\}_{t=1}^\infty$. Then M_T^λ is almost surely well-defined and $\mathbb{E}[M_T^\lambda] \leq 1$.

Proof of Lemma 9. We first show that $\{M_t^\lambda\}_{t=0}^\infty$ is a supermartingale. To see this, we have

$$\begin{aligned}
\mathbb{E}[M_t^\lambda | \mathcal{G}_{t-1}] &= M_{t-1}^\lambda \cdot \mathbb{E} \left[\exp \left(\frac{\langle \lambda^\top X_t, \eta_t \rangle}{\sigma^2 R^2} - \frac{\|\lambda^\top X_t\|_2^2}{2} \right) \middle| \mathcal{G}_{t-1} \right] \\
&\leq M_{t-1}^\lambda \cdot \mathbb{E} \left[\exp \left(\frac{\|\lambda^\top X_t\|_\infty^2}{2R^2} - \frac{\|\lambda^\top X_t\|_\infty^2}{2R^2} \right) \middle| \mathcal{G}_{t-1} \right] \\
&= M_{t-1}^\lambda.
\end{aligned}$$

where the inequality is due to the conditional σ -sub-Gaussianity of η_t and the condition in (44). It also implies $\mathbb{E}[M_t^\lambda] \leq 1$ for any $t \geq 0$. By martingale convergence theorem, with $t \rightarrow \infty$, M_t^λ almost surely converges to a random variable M_∞^λ with finite expectation, and thus M_T^λ is well-defined almost surely. Applying Fatou's lemma, we have

$$\mathbb{E}[M_T^\lambda] = \mathbb{E} \left[\liminf_{t \rightarrow \infty} M_{T \wedge t}^\lambda \right] \leq \liminf_{t \rightarrow \infty} \mathbb{E}[M_{T \wedge t}^\lambda] \leq 1.$$

\square

Lemma 10. *Let T be a stopping time with respect to $\{\mathcal{G}_t\}_{t=0}^\infty$. Then the following holds with probability at least $1 - \delta$*

$$\|S_T\|_{\bar{V}_t}^2 \leq 2\sigma^2 R^2 \log \left(\frac{\det(\bar{V}_t)^{1/2}}{\delta \det(V)^{1/2}} \right).$$

Proof of Lemma 10. For notational simplicity, we assume $\sigma \cdot R = 1$. We define

$$V_t = \sum_{\tau=1}^t \int_{\mathcal{S}} X_{\tau}(s) X_{\tau}(s)^{\top} ds$$

and therefore $\bar{V}_t = V_t + V$. Then we can write $M_t^{\lambda} = \exp(\lambda^{\top} S_t - \|\lambda\|_{V_t}^2/2)$. Let Λ be a \mathbb{R}^d -valued Gaussian random variable with covariance V^{-1} and that it is independent of $\{\mathcal{G}_t\}_{t=0}^{\infty}$. Let $M_t = \mathbb{E}[M_t^{\Lambda} | \mathcal{G}_{\infty}]$ where $\mathcal{G}_{\infty} = \sigma(\cup_{t=0}^{\infty} \mathcal{G}_t)$. Let q denote the density of Λ and $v(A) = \int \exp(-x^{\top} Ax/2) dx = \sqrt{(2\pi)^d / \det(A)}$ for any positive definite matrix A .

Then, we have

$$\begin{aligned} M_t &= \int_{\mathbb{R}^d} \exp(\lambda^{\top} S_t - \|\lambda\|_{V_t}^2/2) q(\lambda) d\lambda \\ &= \int_{\mathbb{R}^d} \exp(-\|\lambda - V_t^{-1} S_t\|_{V_t}^2/2 + \|S_t\|_{V_t^{-1}}^2/2) q(\lambda) d\lambda \\ &= v(V)^{-1} \cdot \exp(\|S_t\|_{V_t^{-1}}^2/2) \cdot \int_{\mathbb{R}^d} \exp\left(-(\|\lambda - V_t^{-1} S_t\|_{V_t}^2 + \|\lambda\|_V^2)/2\right) d\lambda. \end{aligned} \quad (50)$$

Note that

$$\begin{aligned} \|\lambda - V_t^{-1} S_t\|_{V_t}^2 + \|\lambda\|_V^2 &= \|\lambda - \bar{V}_t^{-1} S_t\|_{\bar{V}_t}^2 + \|V_t^{-1} S_t\|_{V_t}^2 - \|S_t\|_{\bar{V}_t^{-1}}^2 \\ &= \|\lambda - \bar{V}_t^{-1} S_t\|_{\bar{V}_t}^2 + \|S_t\|_{V_t^{-1}}^2 - \|S_t\|_{\bar{V}_t^{-1}}^2. \end{aligned} \quad (51)$$

By plugging (51) into (50), we get

$$\begin{aligned} M_t &= v(V)^{-1} \cdot \exp(\|S_t\|_{V_t^{-1}}^2/2) \cdot \int_{\mathbb{R}^d} \exp(-\|\lambda - \bar{V}_t^{-1} S_t\|_{\bar{V}_t}^2/2) d\lambda \\ &= \frac{v(\bar{V}_t)}{v(V)} \cdot \exp(\|S_t\|_{V_t^{-1}}^2/2) \\ &= \sqrt{\det(V)/\det(\bar{V}_t)} \cdot \exp(\|S_t\|_{V_t^{-1}}^2/2). \end{aligned}$$

Hence, we have

$$\mathbb{P}\left(\|S_T\|_{\bar{V}_T^{-1}}^2 > 2 \log\left(\frac{\det(\bar{V}_T)^{1/2}}{\delta \det(V)^{1/2}}\right)\right) = \mathbb{P}(\delta \cdot M_T > 1) \leq \mathbb{E}[\delta \cdot M_T] \leq \delta,$$

which completes the proof of Lemma 10. \square

We construct a stopping time as below.

$$T = \inf \left\{ t \geq 0 : \|S_t\|_{\bar{V}_t^{-1}}^2 > 2 \log\left(\frac{\det(\bar{V}_t)^{1/2}}{\delta \det(V)^{1/2}}\right) \right\}$$

Then, we have

$$\begin{aligned} &\mathbb{P}\left(\exists t \geq 0, \|S_t\|_{\bar{V}_t^{-1}}^2 > 2 \log\left(\frac{\det(\bar{V}_t)^{1/2}}{\delta \det(V)^{1/2}}\right)\right) \\ &= \mathbb{P}(T < \infty) \\ &= \mathbb{P}\left(\|S_T\|_{\bar{V}_T^{-1}}^2 > 2 \log\left(\frac{\det(\bar{V}_T)^{1/2}}{\delta \det(V)^{1/2}}\right), T < \infty\right) \\ &\leq \mathbb{P}\left(\|S_T\|_{\bar{V}_T^{-1}}^2 > 2 \log\left(\frac{\det(\bar{V}_T)^{1/2}}{\delta \det(V)^{1/2}}\right)\right) \\ &\leq \delta, \end{aligned}$$

which completes the proof of Lemma 8. \square

I Experiments

Experiments are conducted on tabular CMOMDPs as follows. We define the constraint set as $\mathcal{W}^* = \{x \in \mathbb{R}^D : \|x\|_2 \leq 1\}$ for simplicity, and one can verify that it satisfies Assumption 1. The transition kernel \mathcal{P} and cost function c are generated uniformly at random from $[0, 1]$ (and we conduct normalization for \mathcal{P}). We make the cost deterministic for simplicity. In addition, we set $\mathcal{P}_h(s_0 | s_0, a_0) = 1$ and $c_h(s_0, a_0) = 0$ for a certain state action pair $(s_0, a_0) \in \mathcal{S} \times \mathcal{A}$ for all $h \in [H]$, and the initial state is set to s_0 . The intuition here is to ensure that the optimal policy, which always takes action a_0 , achieves zero total cost and zero constraint violation for simplicity. The dataset is generated by a uniformly random experimenter, i.e., it picks $a \in \mathcal{A}$ uniformly at random at each step. Hyperparameters are listed in Table 1.

Table 1: List of hyperparameters

Hyperparameter	Value
H : horizon	5
D : dimension of cost function	6
$ \mathcal{S} $: cardinality of state space	5
$ \mathcal{A} $: cardinality of action space	5
$ \mathcal{D} $: dataset size	50000
K : number of iteration of PEDI	100
δ : confidence level	0.9
η : step length	0.01
ν : scaling constant	3

In our implementation, PEDI estimates the transition and cost functions by the empirical mean, i.e., $\hat{\mathcal{P}}_h(s, a) = n_h(s, a, s')/n_h(s, a)$ and $\hat{c}_h^i(s, a) = f_h^i(s, a)/n_h(s, a)$ for $i \in [D]$ where $n_h(s, a)$ is the number of visits to (s, a) at step h and $f_h^i(s, a)$ is the sum of the i -th cost incurred in the dataset when visiting (s, a) at step h . We construct the Hoeffding-style uncertainty quantifiers, i.e., $\Gamma_h^{\mathcal{P}}(s, a, s') = \sqrt{\log(2H|\mathcal{S}||\mathcal{A}||\mathcal{S}|/\delta)/(2n_h(s, a))}$ and $\Gamma_h^c = \sqrt{\log(2DH|\mathcal{S}||\mathcal{A}|/\delta)/(2n_h(s, a))}$. We can verify that they satisfy the definition (Definition 1).

We conduct experiments to see how PEDI converges to the optimal policy with different preference functions: quadratic functions, polynomial functions, and their combinations.

Quadratic Functions. Suppose the interplay of cost functions can be modeled by a positive definite matrix A , a vector b and a constant c , i.e., the preference function is defined as

$$g(x) = \frac{1}{2}x^\top Ax + b^\top x + c,$$

where A is positive definite. For simplicity, we assume b is the zero vector and $c = 0$. To guarantee 1-Lipschitzness, it suffices to restrict the spectral radius λ_{\max} . In particular, we require $\lambda_{\max}(A) \leq 1/(2HD^{1/2})$ since $\|\partial_x g\|_2 = \|2Ax\|_2 \leq 2\lambda_{\max}(A)HD^{1/2}$. For the convex conjugate, we can verify that $g^*(x^*) = \frac{1}{2}(x^* - b)^\top A^{-1}(x^* - b) - c = \frac{1}{2}x^* A^{-1}x^*$. In the numerical experiment, the matrix A is randomly generated with the mentioned spectral radius requirement. The results are given in Table 2.

Table 2: Results of quadratic preference functions

Iteration k	Suboptimality	Constraint Violation
1	0.067	0.880
2	0.505	4.007
3	0.067	0.880
4, 5, ..., 100	0.000	0.000

As we see, it converges to the optimal policy in mere four iterations and stays optimal permanently.

Polynomial Functions. Suppose the preference function is polynomial, i.e.,

$$g(x) = \sum_{i=1}^D c_i |x_i|^{p_i}.$$

For simplicity, we assume $p = p_i = p_j$ and $c = c_i = c_j$ for any $1 \leq i, j \leq D$. To ensure 1-Lipschitzness, it suffices to set $c = 1/(pH^{p-1}D^{1/2})$ for all i which results in $\|\partial_x g\|_2 = cpH^{p-1}D^{1/2} \leq 1$ for $x \geq 0$. Then, we have $g^*(x^*) = \sum_{i=1}^D \frac{|x_i^*|^q}{c^{q-1}p^{q-1}q}$ where $\frac{1}{p} + \frac{1}{q} = 1$. In the numerical experiment, we set $p = 2$. The results are shown in Table 3.

Table 3: Results of polynomial preference functions

Iteration k	Suboptimality	Constraint Violation
1	0.165	1.009
2	0.139	0.844
3, 4, ..., 100	0.000	0.000

As we see, it reaches the optimal solution in only three iterations.

Combination of Quadratic Functions and Polynomial Functions. We consider more complex scenarios where preference functions are combinations of quadratic functions and polynomial functions, i.e., $g(x) = g_1(x_1) + g_2(x_2)$ with

$$g_1(x) = \frac{1}{2}x^\top Ax + b^\top x, \quad g_2(x) = \sum_{i=1}^{D_2} c_i |x_i|^{p_i}.$$

Here $x = (x_1^\top, x_2^\top)^\top$ with $x_1 \in \mathbb{R}^{D_1}$, $x_2 \in \mathbb{R}^{D_2}$ and $D_1 + D_2 = D$. It is clear that $g^*(x^*) = g_1^*(x_1^*) + g_2^*(x_2^*)$. In experiments we set $D_1 = D_2 = D/2$. Moreover, we impose similar requirements and restrictions as we did previously to ensure 1-Lipschitzness. The numerical results are in Table 4, which show that the suboptimality and constraint violation decrease as the number of iterations increases, and PEDI finds the solution in nine iterations.

Table 4: Results of functions that are a combination of quadratic functions and polynomial preference functions

Iteration k	Suboptimality	Constraint Violation
1	1.517	6.033
2	1.490	6.000
3, 4, 5	1.438	5.864
6	1.333	5.609
7	0.811	4.152
8	0.093	0.884
9, 10, ..., 100	0.000	0.000

The following two remarks discuss (1) possibilities to handle other (even more general) preference functions and (2) some practical variants of PEDI for application, which is left to future work, as this paper is mainly theoretical.

Remark 1 (General Preference Function). *In addition to the above demonstration, PEDI is also easily applicable to preference functions from many function classes such as exponential function, logarithmic function, and entropy function. Even when the exact expression of the preference function g is not good or even unknown, PEDI applies as long as we can approximate g^* by some numerical methods, say, by directly approximating $\sup_x (\langle x^*, x \rangle - g(x))$, which is the definition of the convex conjugate. To obtain the subgradient, we can use certain techniques such as numerical differentiation.*

Remark 2 (General Planning Algorithm). *For a real-world application, the pessimistic planning (PESSPLANNING, see Algorithm 1) may seem too heavy. It can be replaced by any algorithms as long as it approximately produces the desired policy π^k and a pessimistic estimation of the value functions V^k at each iteration. For example, we can apply policy iteration algorithms or even any neural network-based approximate algorithms.*