

---

# On Dynamic Programming Decompositions of Static Risk Measures in Markov Decision Processes

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Optimizing static risk-averse objectives in Markov decision processes is difficult  
2 because they do not admit standard dynamic programming equations common in  
3 Reinforcement Learning (RL) algorithms. Dynamic programming decompositions  
4 that augment the state space with discrete risk levels have recently gained pop-  
5 ularity in the RL community. Prior work has shown that these decompositions  
6 are optimal when the risk level is discretized sufficiently. However, we show that  
7 these popular decompositions for Conditional-Value-at-Risk (CVaR) and Entropic-  
8 Value-at-Risk (EVaR) are inherently suboptimal regardless of the discretization  
9 level. In particular, we show that a saddle point property assumed to hold in prior  
10 literature may be violated. However, a decomposition does hold for Value-at-Risk  
11 and our proof demonstrates how this risk measure differs from CVaR and EVaR.  
12 Our findings are significant because risk-averse algorithms are used in high-stake  
13 environments, making their correctness much more critical.

## 14 1 Introduction

15 Risk-averse reinforcement learning (RL) seeks to provide a risk-averse policy for high stake real-  
16 world decision problems. These high-stake domains include autonomous driving (Jin et al., 2019;  
17 Sharma et al., 2020), robot collision avoidance (Ahmadi et al., 2021; Hakobyan and Yang, 2021),  
18 liver transplant timing (Köse, 2016), HIV treatment (Keramati et al., 2020; Zhong, 2020), unmanned  
19 aerial vehicle (UAV) (Choudhry et al., 2021), and investment liquidation (Min et al., 2022), to name a  
20 few. Because these domains call for reliable solutions, risk-averse algorithms must be based on solid  
21 theoretical foundations. This is one reason why monetary risk measures, such as Value-at-Risk (VaR)  
22 and Conditional Value-at-Risk (CVaR), have become pervasive in risk-averse RL (Prashanth and  
23 Fu, 2022). Indeed, risk measures such as CVaR are known to be coherent (Artzner et al., 1999)  
24 with respect to a set of fundamental axioms that define how risk should be quantified and have been  
25 adopted as gold standards in banking regulations (Basel Committee on Banking Supervision, 2019).

26 Introducing risk-averse objectives in Markov decision processes (MDPs), the primary model used  
27 in RL, is challenging. Dynamic programming, the linchpin of most RL algorithms, cannot be used  
28 directly to optimize a risk measure like VaR or CVaR in MDPs. One line of work tackles this  
29 challenge by exploiting the primal representation of risk measures and augmenting the state space  
30 of their dynamic programs (DPs) with an additional parameter that typically represents the total  
31 cumulative reward up to the current point (Bäuerle and Ott, 2011; Boda et al., 2004; Chow and  
32 Ghavamzadeh, 2014; Filar et al., 1995; Hau et al., 2023; Lin et al., 2003; Wu and Lin, 1999; Xu and  
33 Mannor, 2011). Even when the original MDP is finite, this DP requires computing the value function  
34 for a continuous state space, and thus, has been considered inefficient in practice (Chapman et al.,  
35 2022; Chow et al., 2015; Li et al., 2022).

36 Another line of recent work leverages the dual representation to produce a *risk-level decomposition*  
37 of risk measures (Pflug and Pichler, 2016). Using this decomposition, numerous authors have derived

38 DPs for common risk measures and integrated them within various RL algorithms (Chapman et al.,  
39 2019, 2022; Chow et al., 2015; Ding and Feinberg, 2022; Ni and Lai, 2022; Rigter et al., 2021; Stanko  
40 and Macek, 2019). Although this risk-level decomposition requires augmenting the state space with a  
41 continuous parameter, this parameter is naturally bounded between 0 and 1. It has been generally  
42 accepted, with several *tentative* proofs supporting this claim (Chow et al., 2015; Li et al., 2022),  
43 that these DPs recover the optimal policy if we can discretize the augmented state space sufficiently  
44 finely. Moreover, it is believed that one can use the optimal value function from this DP to recover  
45 the policies that are optimal for the full range of risk levels.

46 In this paper, we make a surprising discovery that numerous claims of optimality of risk-level  
47 decompositions published in the past several years are incorrect. Even when one discretizes the  
48 augmented state space arbitrarily finely, most risk-level DPs are not guaranteed to recover the optimal  
49 value function and policy. There are several reasons why existing arguments fail. As the most common  
50 reason, several papers assume that a certain saddle point property holds, either explicitly (Chow  
51 et al., 2015) or implicitly (Ding and Feinberg, 2022; Li et al., 2022). We show that this property does  
52 not generally hold, invalidating the optimality of DPs, as hinted at in Chapman et al. (2019, 2022).  
53 This finding directly refutes the claimed or hypothesized optimality of algorithms proposed in many  
54 recent research papers and pre-prints, such as Chapman et al. (2019); Chow et al. (2015); Ding and  
55 Feinberg (2022); Li et al. (2022); Rigter et al. (2021); Stanko and Macek (2019). Our results also  
56 affect applications of these algorithms, such as automated vehicle motion planning (Jin et al., 2019).  
57 We also identify gaps in related decompositions (Li et al., 2022; Ni and Lai, 2022) and propose how  
58 to fix them.

59 We make the following contributions in this paper. *First*, we show in Section 3 that the popular DP  
60 for optimizing CVaR in MDPs may not recover the optimal value function and policy regardless of  
61 how finely one discretizes the risk level in the augmented states. This method was first proposed  
62 in Chow et al. (2015) but adopted widely afterwards (Chapman et al., 2019; Ding and Feinberg,  
63 2022; Li et al., 2022; Rigter et al., 2021; Stanko and Macek, 2019). The simple counterexample in  
64 this section contradicts the optimality claims in Chow et al. (2015); Ding and Feinberg (2022); Li  
65 et al. (2022). We hypothesize that prior work missed this issue, because the CVaR DP works for  
66 policy evaluation and only fails when one uses it to optimize policies based on the “risk-to-go value  
67 function”. Therefore, our results do not contradict the original decomposition in Pflug and Pichler  
68 (2016) that only applies to policy evaluation. We give a new independent and simple proof that the  
69 CVaR decomposition indeed works when evaluating a fixed policy.

70 *Second*, we show in Section 4 that the DP for optimizing the Entropic-Value-at-Risk in MDPs,  
71 proposed by Ni and Lai (2022), does not compute the correct value function even when the policy  
72 is fixed. Although EVaR has not been as popular as CVaR, it has been gaining attention in recent  
73 years (Hau et al., 2023). We give an example that contradicts the correctness claims of the risk-level  
74 decomposition for EVaR in Ni and Lai (2022). The gap that we identify with this objective applies  
75 to both policy evaluation and policy optimization. Furthermore, we prove a new, correct EVaR  
76 decomposition for policy evaluation. Unfortunately, the EVaR decomposition fails and is sub-optimal  
77 when is applied to policy optimization, similar to CVaR.

78 *Third*, we propose an *optimal dynamic program* for policy optimization of VaR in Section 5. Our DP is  
79 based on a risk-level decomposition that closely resembles the quantile MDP decomposition in Li et al.  
80 (2022) but corrects for several technical inaccuracies. The derivation shows why VaR stands apart  
81 from coherent risk measures like CVaR and EVaR. VaR is unique in that the decomposition can be  
82 constructed directly from the primal formulation of the risk measure, which avoids the complications  
83 that arise in the robust formulations used in CVaR and EVaR decompositions.

84 It is important to note that the correctness of DPs that augment the state space with the accumulated  
85 rewards is unaffected by our results (Bäuerle and Ott, 2011; Chow and Ghavamzadeh, 2014; Chow  
86 et al., 2018; Hau et al., 2023). These DPs use the *primal* risk measure representation and do not suffer  
87 from the same saddle point issue as the augmentation methods that use the *dual* representation of the  
88 risk measures, such as the one in Chow et al. (2015).

## 89 2 Preliminaries

90 This section summarizes relevant properties of monetary risk measures and outlines how they are  
91 typically used in the context of solving MDPs.

92 **Monetary Risk Measures** We restrict our attention to probability spaces with a finite outcome  
 93 space  $\Omega$  such that  $|\Omega| = m$  for some  $m \in \mathbb{N}$ . We use  $\mathbb{X} = \mathbb{R}^m$  to denote the space of real-valued  
 94 random variables. To improve the clarity of probabilistic claims, we always adorn random variables  
 95 with a tilde, such as  $\tilde{x} \in \mathbb{X}$ . In finite spaces, we can represent any random variable  $\tilde{x} \in \mathbb{X}$  as a vector  
 96  $\mathbf{x} \in \mathbb{R}^m$ . We also use  $\mathbf{q} \in \Delta_m$  to represent a probability distribution over  $\Omega$  where  $\Delta_m$  represents  
 97 the  $m$ -dimensional probability simplex. Using this notation, we can write that  $\mathbb{E}[\tilde{x}] = \mathbf{q}^\top \mathbf{x}$ .

98 A *monetary risk measure*  $\psi: \mathbb{X} \rightarrow \mathbb{R}$  assigns a real value to each real-valued random variable in a  
 99 way that it is monotone and cash-invariant (Follmer and Schied, 2016; Shapiro et al., 2014). A risk  
 100 measure can be seen as a generalization of the expectation operator  $\mathbb{E}[\cdot]$  that also takes into account  
 101 the uncertainty in the random variable. In this work, we define all risk measures for random variables  
 102  $\tilde{x}$  that represent *rewards*. Thus, the risk-averse decision-maker aims to choose actions that maximize  
 103 the value of the risk measure, i.e., a higher value of risk measure represents a lower exposure to risk.

104 We consider three monetary risk measures common in RL. Perhaps the most well-known measure  
 105 is *Value-at-Risk* (VaR), which is defined for a risk-level  $\alpha \in [0, 1]$  and a random variable  $\tilde{x} \in \mathbb{X}$  in  
 106 modern literature as (e.g., Follmer and Schied 2016; Shapiro et al. 2014)

$$\text{VaR}_\alpha[\tilde{x}] = \sup \{z \in \mathbb{R} \mid \mathbb{P}[\tilde{x} < z] \leq \alpha\} = \inf \{z \in \mathbb{R} \mid \mathbb{P}[\tilde{x} \leq z] > \alpha\}. \quad (1)$$

107 Note that  $\text{VaR}_1[\tilde{x}] = \infty$ . The equality between the two definition holds, for example, by Follmer and  
 108 Schied (2016, remark A.20). Another popular risk measure is the *Conditional-value-at-Risk* (CVaR),  
 109 which is defined for a risk level  $\alpha \in [0, 1]$  and a random variable  $\tilde{x} \in \mathbb{X}$  distributed as  $\tilde{x} \sim \mathbf{q}$   
 110 as (e.g., Follmer and Schied 2016, definition 11.8 and Shapiro et al. 2014, eq. 6.23)

$$\text{CVaR}_\alpha[\tilde{x}] = \sup_{z \in \mathbb{R}} (z - \alpha^{-1} \mathbb{E}[z - \tilde{x}]_+) = \inf_{\boldsymbol{\xi} \in \Delta_m} \left\{ \boldsymbol{\xi}^\top \mathbf{x} \mid \alpha \cdot \boldsymbol{\xi} \leq \mathbf{q} \right\}, \quad (2)$$

111 with  $\text{CVaR}_0[\tilde{x}] = \text{ess inf}[\tilde{x}]$  and  $\text{CVaR}_1[\tilde{x}] = \mathbb{E}[\tilde{x}]$ . Finally, the *entropic value at risk* (EVaR), with  
 112  $\text{EVaR}_0[\tilde{x}] = \text{ess inf}[\tilde{x}]$  and  $\text{EVaR}_1[\tilde{x}] = \mathbb{E}[\tilde{x}]$ , is defined for  $\alpha \in (0, 1]$  as (Ahmadi-Javid, 2012)

$$\text{EVaR}_\alpha[\tilde{x}] = \sup_{\beta > 0} \frac{1}{\beta} \left( -\log \alpha^{-1} \mathbb{E}[\exp(-\beta \tilde{x})] \right) = \inf_{\boldsymbol{\xi} \in \Delta_m: \boldsymbol{\xi} \ll \mathbf{q}} \left\{ \boldsymbol{\xi}^\top \mathbf{x} \mid \text{KL}(\boldsymbol{\xi} \parallel \mathbf{q}) \leq -\log \alpha \right\}, \quad (3)$$

113 where KL is the standard KL-divergence defined for each  $\mathbf{x}, \mathbf{y} \in \Delta_m$  as  $\text{KL}(\mathbf{x} \parallel \mathbf{y}) =$   
 114  $\sum_{\omega \in \Omega} x_\omega \log(x_\omega/y_\omega)$ . This definition is valid only when  $\mathbf{x}$  is absolutely continuous with respect to  
 115  $\mathbf{y}$ , which is denoted as  $\mathbf{x} \ll \mathbf{y}$  and corresponds to  $y_\omega = 0 \Rightarrow x_\omega = 0$  for each  $\omega \in \Omega$ .

116 **Risk Averse MDPs** A Markov decision process (MDP) is a sequential decision model that underlies  
 117 most of RL (Puterman, 2005). We consider finite MDPs with states  $\mathcal{S} = \{s_1, \dots, s_S\}$  and actions  
 118  $\mathcal{A} = \{a_1, \dots, a_A\}$ . After taking an action in a state, the agent transitions to a next state according  
 119 to a transition probability function  $p: \mathcal{S} \times \mathcal{A} \rightarrow \Delta_{\mathcal{S}}$  such that  $p(s, a, s')$  represents the transition  
 120 probability from  $s \in \mathcal{S}$  to  $s' \in \mathcal{S}$  after taking  $a \in \mathcal{A}$ . We use  $\mathbf{p}_{s,a} = p(s, a, \cdot) \in \Delta_{\mathcal{S}}$  to denote the  
 121 vector of transition probabilities. The initial state  $\tilde{s}_0$  is distributed according to  $\hat{\mathbf{p}} \in \Delta_{\mathcal{S}}$ . To avoid  
 122 divisions by 0 that are not central to our claims, we assume that  $\hat{p}_s > 0$  for each  $s \in \mathcal{S}$ . Finally,  
 123 the reward function is  $r: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ , where  $r(s, a, s')$  represents the deterministic reward  
 124 associated with the transition to  $s'$  from  $s$  after taking an action  $a$ .

125 The most-general solution to an MDP is a *history-dependent randomized* policy  $\pi$  which maps a  
 126 sequence of observed states and actions  $s^0, a^0, s^1, a^1, \dots, s^t$  to a distribution over the next action  $a^t$ .  
 127 It is well-known that with risk-neutral objectives, there always exists an optimal stationary (depends  
 128 only on the last state) deterministic policy (Puterman, 2005). When the objective is risk-averse, like  
 129 VaR, or CVaR, there may not exist an optimal stationary or deterministic policy. Hence, we use the  
 130 symbol  $\Pi$  to denote the set of history-dependent randomized policies in the remainder of the paper.

131 This paper focuses on the *finite-horizon objective* in which the agent aims to compute policies  
 132 that optimize the sum of rewards over a known horizon  $T$ . We further restrict our attention to the  
 133 objective with horizon  $T = 1$ . It turns out that having a single time-step is sufficient to derive our  
 134 counterexamples to existing dynamic programs. Moreover, deriving the decompositions with  $T = 1$   
 135 makes it possible to avoid technicalities caused by history-dependent policies, which could distract  
 136 us from the main ideas presented in this work. Our results can be extended to general horizons  $T > 1$   
 137 and the discounted infinite-horizon objectives using standard techniques (Chow et al., 2015).

138 With horizon  $T = 1$ , the set of randomized history-dependent policies is  $\Pi = \{\boldsymbol{\pi}: \mathcal{S} \rightarrow \Delta_{\mathcal{A}}\}$ . The  
 139 symbol  $\pi(s, a)$  denotes the probability of an action  $a$  in a state  $s$ , and  $\boldsymbol{\pi}(s) = \pi(s, \cdot) \in \Delta_{\mathcal{A}}$  denotes

140 the  $A$ -dimensional vector of action probabilities in a state  $s$ . Given a risk measure  $\psi$  with a risk level  
 141  $\alpha \in [0, 1]$ , the finite-horizon risk-averse value of a policy  $\pi \in \Pi$  is computed as

$$v_0^\pi(\alpha) := \psi_\alpha^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}')], \quad (4)$$

142 where the superscript in  $\psi_\alpha^{\tilde{a} \sim \pi(\tilde{s})}$  specifies the distribution of the random action. Throughout the  
 143 paper, we generally use  $\tilde{s}$  to denote the random state at time  $t = 0$  and  $\tilde{s}'$  to denote the random state at  
 144 time  $t = 1$ . In risk-neutral objectives, when  $\psi = \mathbb{E}$ , one can use the tower property of the expectation  
 145 operator and define a value function  $v_t$  for each time step  $t$  (Puterman, 2005), but this property does  
 146 not hold in most static risk measures (Hau et al., 2023). The term *policy evaluation* in the remainder  
 147 of the paper refers to computing the value in (4).

148 The goal in an MDP is to compute an *optimal* value function and a policy that attains it. In risk-averse  
 149 MDPs, this goal is formalized as the following risk-averse optimization

$$v_0^*(\alpha) := \max_{\pi \in \Pi} v_0^\pi(\alpha) = \max_{\pi \in \Pi} \psi_\alpha^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}')], \quad (5)$$

150 with the *optimal policy*  $\pi^*$  being any policy that attains the maximum in (5). As with policy evaluation,  
 151 when  $\psi = \mathbb{E}$ , the optimal value function  $v_t^*$  can be defined for each time-step  $t$  (Puterman, 2005),  
 152 but this is impossible in general for common risk measures, like VaR and CVaR. The term *policy*  
 153 *optimization* in the remainder of the paper refers to computing the value and the maximizer in (5).

154 In the remainder of the paper, we study dynamic programming algorithms proposed to solve the  
 155 policy evaluation problem in (4) and policy optimization problem in (5). In general, these algorithms  
 156 build on risk-level decomposition (Pflug and Pichler, 2016) of risk measures to define a value function  
 157  $v_t^\pi(s, \alpha)$  for each time step  $t \in [T]$ , state  $s \in \mathcal{S}$ , and risk-level  $\alpha \in [0, 1]$  (Chow et al., 2015). The  
 158 value function represents the risk-adjusted sum of rewards that can be obtained if starting in a state  
 159  $s \in \mathcal{S}$  at time  $t$  and a risk level  $\alpha$ . For example, one would define the value function as  $v_1^\pi(s, \alpha) :=$   
 160  $\psi_\alpha^{\tilde{a} \sim \pi(s)} [r(s, \tilde{a}, \tilde{s}')]$  and compute  $v_0^\pi$  using a Bellman operator  $T_\alpha^\pi$  as  $v_0^\pi(\alpha) = (T_\alpha^\pi v_1^\pi)(\alpha)$ . In risk-  
 161 neutral objectives, the Bellman operator is defined as  $(T^\pi(v_1^\pi))(\alpha) = \mathbb{E}[v_1^\pi(\tilde{s}, \alpha)]$ , with  $\alpha \in \{1\}$ ,  
 162 but in risk-averse formulations the operator definition is more complex. The remainder of the paper  
 163 discusses the decompositions and the operator for CVaR, EVaR, and VaR risk measures respectively.

### 164 3 CVaR: Decomposition Fails in Policy Optimization

165 In this section, we show that a common CVaR decomposition proposed in Chow et al. (2015) and used  
 166 to optimize risk-averse policies is inherently sub-optimal regardless of how closely one discretizes  
 167 the state space. The following proposition represents one of the key results used to decompose the  
 168 risk measure in multi-stage decision-making.

169 **Proposition 3.1** (lemma 22 in Pflug and Pichler 2016). *Suppose that  $\pi \in \Pi$  and  $\tilde{s} \sim \hat{p}$ ,  $\tilde{a} \sim \pi(\tilde{s})$ ,*  
 170  *$\tilde{s}' \sim \mathbf{p}_{s,\alpha}$ . Then,*

$$\text{CVaR}_\alpha [r(\tilde{s}, \tilde{a}, \tilde{s}')] = \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s], \quad (6)$$

171 where the state  $s$  on the right-hand side is not random and

$$\mathcal{Z}_C = \{\zeta \in \Delta_S \mid \alpha \cdot \zeta \leq \hat{p}\}. \quad (7)$$

172 The notation in Proposition 3.1 differs superficially from lemma 22 in Pflug and Pichler (2016).  
 173 Specifically, our CVaR is defined for rewards rather than costs, the meaning of our  $\alpha$  corresponds to  
 174  $1 - \alpha$  in Pflug and Pichler (2016), and we use  $\xi_s = z_s \hat{p}_s$  as the optimization variable. We include a  
 175 simple proof of Proposition 3.1 for completeness in Appendix A.1.

176 The decomposition in Proposition 3.1 is important because it shows that the CVaR evaluation can be  
 177 formulated as a dynamic program. The theorem shows that CVaR at time  $t = 0$  decomposes into  
 178 a convex combination of CVaR values at time  $t = 1$ . Recursively repeating this process, one can  
 179 formulate a dynamic program for any finite time horizon  $T$ . Because the risk-level at time  $t = 1$   
 180 differs from the level at  $t = 0$  and depends on the optimal  $\zeta$ , one must compute CVaR values for  
 181 all (or many) risk-levels  $\alpha \in [0, 1]$  at time  $t = 1$ . As a result, the dynamic program includes an  
 182 additional state variable that represents the current risk level.

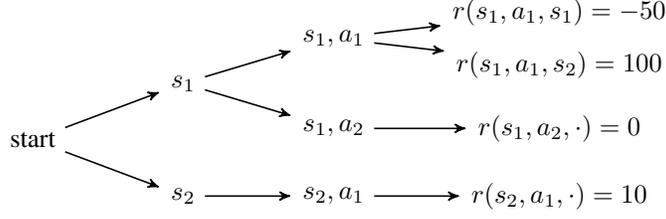


Figure 1: Rewards of MDP  $M_C$  used in the proof of Theorem 3.2. The dot indicates that the rewards are independent of the next state.

183 [Chow et al. \(2015\)](#) proposed to adapt the decomposition in Proposition 3.1 to policy optimization as

$$\begin{aligned} \max_{\pi \in \Pi} \text{CVaR}_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}')] &= \max_{\pi \in \Pi} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \left( \text{CVaR}_{\alpha \xi_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(s)} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \right) \\ &\stackrel{??}{=} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \left( \max_{d \in \Delta_{\mathcal{A}}} \text{CVaR}_{\alpha \xi_s \hat{p}_s^{-1}}^{\tilde{a} \sim d} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \right). \end{aligned} \quad (8)$$

184 They used the decomposition in (8) to formulate a dynamic program with the current risk-level as an  
 185 additional state variable. We prove in the following theorem that the second equality in (8) marked  
 186 with question marks is false in general.

187 **Theorem 3.2.** *There exists an MDP and a risk level  $\alpha \in [0, 1]$  such that*

$$\max_{\pi \in \Pi} \text{CVaR}_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}')] < \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \left( \max_{d \in \Delta_{\mathcal{A}}} \text{CVaR}_{\alpha \xi_s \hat{p}_s^{-1}}^{\tilde{a} \sim d} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \right). \quad (9)$$

188 Before proving Theorem 3.2, we discuss its implications. First, Theorem 3.2 contradicts theorems 5  
 189 and 7 in [Chow et al. \(2015\)](#) and shows that their algorithm is inherently sub-optimal regardless of  
 190 the resolution of the discretization. Theorem 3.2 also contradicts the optimality of the accelerated  
 191 dynamic program proposed in [Stanko and Macek \(2019\)](#). The result of [Chow et al., 2015](#)  
 192 was exploited as is in [Chapman et al. \(2019\)](#), [Ding and Feinberg \(2022\)](#), and [Jin et al. \(2019\)](#) to propose  
 193 DP reductions, and extended, without proof, in [Rigter et al., 2021](#) to the context of a Bayesian MDP.

194 Finally, it is important to emphasize that Theorem 3.2 only applies to the policy optimization setting  
 195 and does not contradict Proposition 3.1, which holds for the evaluation of policies that assign the  
 196 same action distribution to each history of states and actions (i.e., policies that are independent of the  
 197 hypothesized values of  $\zeta$ ).

198 *Proof.* Let  $\alpha = 0.5$  and consider the MDP  $M_C$  in Figure 1. In state  $s_1$ , both actions  $a_1$  and  $a_2$  are  
 199 available, and in state  $s_2$ , only action  $a_1$  is available. The MDP's rewards are

$$\begin{aligned} r(s_1, a_1, s_1) &= -50, & r(s_1, a_1, s_2) &= 100, \\ r(s_1, a_2, s_1) &= r(s_1, a_2, s_2) = 0, & r(s_2, a_1, s_1) &= r(s_2, a_1, s_2) = 10. \end{aligned}$$

200 The transition probabilities in  $M_C$  are

$$p(s_1, a_1, s_1) = 0.4, \quad p(s_1, a_1, s_2) = 0.6,$$

201 and the initial distribution is uniform:  $\hat{p}_{s_1} = \hat{p}_{s_2} = 0.5$ .

202 To simplify the notation, we define  $\theta_{\pi}: \mathcal{Z}_C \rightarrow \mathbb{R}$  for each  $\pi \in \Pi$  and  $\zeta \in \mathcal{Z}_C$  as

$$\theta_{\pi}(\zeta) = \sum_{s \in \mathcal{S}} \zeta_s \text{CVaR}_{\alpha \xi_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(s)} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s].$$

203 Because CVaR is convex in the distribution ([Delage et al., 2019](#)) and any distribution for  $r(\tilde{s}, \tilde{a}, \tilde{s}')$   
 204 obtained from a policy  $\pi \in \Pi$  is a mixture of the distributions of  $r(\tilde{s}, a_1, \tilde{s}')$  and  $r(\tilde{s}, a_2, \tilde{s}')$ , it is  
 205 sufficient to consider only *deterministic* policies (there exists an optimal deterministic policy). Thus,

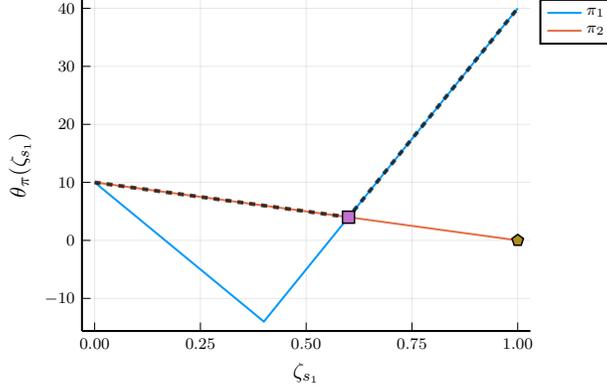


Figure 2: The functions  $\theta_{\pi_1}(\cdot)$  and  $\theta_{\pi_2}(\cdot)$  used in the CVaR counterexample in the proof of Theorem 3.2. The dashed line shows the function  $\zeta_{s_1} \mapsto \max_{\pi \in \{\pi_1, \pi_2\}} \theta_{\pi}(\zeta_{s_1}, 1 - \zeta_{s_1})$ .

206 we can reformulate the *left-hand side* of (9) in terms of  $\theta_{\pi}(\zeta)$  as

$$\begin{aligned} \max_{\pi \in \Pi} \text{CVaR}_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}')] &= \max_{\pi \in \{\pi_1, \pi_2\}} \text{CVaR}_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}')] \\ &= \max_{\pi \in \{\pi_1, \pi_2\}} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \cdot \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(s)} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \\ &= \max_{\pi \in \{\pi_1, \pi_2\}} \min_{\zeta \in \mathcal{Z}_C} \theta_{\pi}(\zeta), \end{aligned}$$

207 with  $\pi_1(s, a_1) = 1 - \pi_1(s, a_2) = 1$  and  $\pi_2(s, a_2) = 1 - \pi_2(s, a_1) = 1$ , for all  $s \in \mathcal{S}$ . The functions  
208  $\theta_{\pi_1}(\cdot)$  and  $\theta_{\pi_2}(\cdot)$  are depicted in Figure 2. Similarly, the *right-hand side* of (9) can be expressed  
209 using the convexity of CVaR in the distribution by algebraic manipulation as

$$\min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \max_{d \in \Delta_A} \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim d} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] = \min_{\zeta \in \mathcal{Z}_C} \max_{\pi \in \{\pi_1, \pi_2\}} \theta_{\pi}(\zeta).$$

210 Using the notation introduced above and the sufficiency of optimizing over deterministic policies  
211 only, the inequality in (9) becomes

$$\max_{\pi \in \{\pi_1, \pi_2\}} \min_{\zeta \in \mathcal{Z}_C} \theta_{\pi}(\zeta) < \min_{\zeta \in \mathcal{Z}_C} \max_{\pi \in \{\pi_1, \pi_2\}} \theta_{\pi}(\zeta). \quad (10)$$

212 Figure 2 demonstrates the inequality in (10) numerically, with the rectangle representing the left-  
213 hand side maximum and the pentagon representing the right-hand side minimum. The dashed line  
214 represents the function  $\zeta \mapsto \max_{\pi \in \{\pi_1, \pi_2\}} \theta_{\pi}(\zeta)$ .

215 To show the strict inequality in (10) formally, we evaluate the functions  $\theta_{\pi_1}(\cdot)$  and  $\theta_{\pi_2}(\cdot)$  for MDP  
216  $M_C$ . The function  $\theta_{\pi_2}(\cdot)$  is linear because the CVaR applies to a constant, and CVaR is translation  
217 invariant. The function  $\theta_{\pi_1}(\cdot)$  is piecewise-linear and convex, and its slope can be computed using  
218 the subgradient that for each  $s \in \mathcal{S}$  and  $\hat{\zeta} \in \mathcal{Z}_C$  satisfies (Chow et al., 2015)

$$\partial_{\zeta_s} \hat{\zeta}_s \text{CVaR}_{\alpha \hat{p}_s^{-1} \hat{\zeta}_s} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \ni \text{VaR}_{\alpha \hat{p}_s^{-1} \hat{\zeta}_s} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s].$$

219 Simple algebraic manipulation then shows that

$$\theta_{\pi_1}(\zeta) = \max \{10 - 60 \zeta_{s_1}, 90 \zeta_{s_1} - 50\}, \quad \theta_{\pi_2}(\zeta) = 10 - 10 \zeta_{s_1},$$

220 and  $\mathcal{Z}_C = \Delta_{\mathcal{S}}$ , which implies that  $\zeta_{s_1} \in (0, 1)$ . Therefore, by algebraic manipulation, we get the  
221 desired strict inequality

$$0 = \max_{\pi \in \{\pi_1, \pi_2\}} \min_{\zeta \in \mathcal{Z}_C} \theta_{\pi}(\zeta) < \min_{\zeta \in \mathcal{Z}_C} \max_{\pi \in \{\pi_1, \pi_2\}} \theta_{\pi}(\zeta) = 4,$$

222 where 0 and 4 are represented by the pentagon and rectangle in Figure 2, respectively.  $\square$

223 In summary, the decomposition in Proposition 3.1 cannot be exploited in policy optimization because  
 224 the inequality in the derivation above may not be tight:

$$\begin{aligned} \max_{\pi \in \Pi} \text{CVaR}_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})} [r(\tilde{s}, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] &= \max_{\pi \in \Pi} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(\tilde{s})} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \\ &\leq \min_{\zeta \in \mathcal{Z}_C} \max_{\pi \in \Pi} \sum_{s \in \mathcal{S}} \zeta_s \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(\tilde{s})} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \\ &= \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \max_{d \in \Delta_A} \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim d} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s], \end{aligned}$$

225 where the last equality follows from the interchangeability property of optimization and expected  
 226 value (Shapiro et al., 2014, Theorem 7.92).

227 It is finally worth noting that we omit to comment on the validity of the CVaR decomposition in Li et al.  
 228 (2022) given that it considers a different measure than the CVaR defined in Equation (1). Namely their  
 229 measure takes the form:  $\widetilde{\text{CVaR}}_{\alpha}[\tilde{x}] := \inf_{z \in \mathbb{R}} (z + (1 - \alpha)^{-1} \mathbb{E}[\tilde{x} - z]_+) = -\text{CVaR}_{1-\alpha}[-\tilde{x}]$ ,  
 230 which is not a coherent risk measure.

## 231 4 EVaR: Decomposition Fails for Policy Evaluation

232 In this section, we show that a decomposition for EVaR proposed in Ni and Lai (2022) is inexact even  
 233 when considering the policy evaluation setting. Ni and Lai (2022) recently proposed a decomposition  
 234 of EVaR for a fixed  $\pi \in \Pi$  with  $\tilde{a} \sim \pi(\tilde{s})$  and a risk level  $\alpha \in (0, 1]$  as

$$\text{EVaR}_{\alpha} [r(\tilde{s}, \tilde{a}, \tilde{s}')] \stackrel{??}{=} \min_{\xi \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s], \quad (11)$$

235 where  $\tilde{s} \sim \hat{p}$ ,  $\tilde{s}' \sim p(\tilde{s}, \tilde{a}, \cdot)$ , and

$$\mathcal{Z}_E = \left\{ \xi \in \Delta_{\mathcal{S}} \mid \sum_{s \in \mathcal{S}} \xi_s \log(\xi_s / \hat{p}_s) \leq -\log \alpha, \quad \overbrace{\alpha \cdot \xi \leq \hat{p}}^{\text{implicit in Ni and Lai (2022)}} \right\}. \quad (12)$$

236 Note that we use variables  $\xi_s = z_s \hat{p}_s$  in comparison with  $z_s$  in Ni and Lai (2022).

237 The constraint  $\alpha \cdot \xi \leq \hat{p}$  in (12) was not stated explicitly in Ni and Lai (2022) but is necessary  
 238 because  $\text{EVaR}_{\alpha'}[\cdot]$  is defined only for  $\alpha' \in [0, 1]$ . When  $\alpha' = \alpha \xi_s \hat{p}_s^{-1}$  in (11) it must also satisfy  
 239 for each  $s \in \mathcal{S}$  that

$$\alpha' \leq 1 \Leftrightarrow \alpha \xi_s \hat{p}_s^{-1} \leq 1 \Leftrightarrow \alpha \cdot \xi_s \leq \hat{p}_s.$$

240 This additional constraint on  $\xi$  implies that  $\mathcal{Z}_E \subseteq \mathcal{Z}_C$ , for the  $\mathcal{Z}_C$  defined in (7).

241 We claim in the following theorem (see Appendix A.2 for a proof) that the equality in (11) does not  
 242 hold even in the policy evaluation setting.

243 **Theorem 4.1.** *There exists an MDP with a single action and  $\alpha \in (0, 1]$  such that*

$$\text{EVaR}_{\alpha} [r(\tilde{s}, a_1, \tilde{s}')] < \min_{\xi \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}} [r(s, a_1, \tilde{s}') \mid \tilde{s} = s], \quad (13)$$

244 *the set  $\mathcal{Z}_E$  defined by (12).*

245 Theorem 4.1 demonstrates a stronger failure mode than the one in Theorem 3.2 (for CVaR policy  
 246 optimization), since it applies to both policy evaluation and policy optimization settings.

247 We propose a correct decomposition of EVaR in the following theorem and employ it to establish that  
 248 the decomposition in (11) overestimates the actual value of EVaR (see Appendix A.3 for a proof).

249 **Theorem 4.2.** *Given any finite MDP with horizon  $T = 1$  and  $\alpha \in (0, 1]$ , we have that*

$$\text{EVaR}_{\alpha} [r(\tilde{s}, \tilde{a}, \tilde{s}')] = \inf_{\zeta \in (0, 1]^{\mathcal{S}}, \xi \in \mathcal{Z}'_E(\zeta)} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\zeta_s} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s],$$

250 *where*

$$\mathcal{Z}'_E(\zeta) = \left\{ \xi \in \Delta_{\mathcal{S}} \mid \xi \ll \hat{p}, \sum_{s \in \mathcal{S}} \xi_s (\log(\xi_s / \hat{p}_s) - \log(\zeta_s)) \leq -\log \alpha \right\}.$$

251 *Moreover, EVaR can be upper-bounded as*

$$\text{EVaR}_{\alpha} [r(\tilde{s}, \tilde{a}, \tilde{s}')] \leq \min_{\xi \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s]. \quad (14)$$

## 252 5 VaR: Decomposition Holds for Policy Evaluation and Optimization

253 In this section, we discuss a dynamic program decomposition for VaR whose decomposition resembles  
 254 those for CVaR and EVaR described in Sections 3 and 4. We provide a new proof of the VaR  
 255 decomposition to elucidate the differences that make it optimal in contrast to CVaR and EVaR  
 256 decompositions. Our VaR decomposition closely resembles the quantile MDP approach in Li et al.  
 257 (2022) with a few technical modifications that can significantly impact the computed value.

258 To contrast the typical definition of VaR with the quantile definition in Li et al. (2022), it is helpful  
 259 to summarize how VaR is related to the quantile of a random variable. Let  $q \in \mathbb{R}$  define as the  
 260  $\alpha$ -quantile of  $\tilde{x} \in \mathbb{X}$  when

$$\mathbb{P}[\tilde{x} \leq q] \geq \alpha \quad \text{and} \quad \mathbb{P}[\tilde{x} < q] \leq \alpha. \quad (15)$$

261 In general, the set of quantiles is an interval  $[q_{\tilde{x}}^-(\alpha), q_{\tilde{x}}^+(\alpha)]$  with the bounds computed as (Follmer  
 262 and Schied, 2016, appendix A.3)

$$\begin{aligned} q_{\tilde{x}}^-(\alpha) &= \sup \{z \mid \mathbb{P}[\tilde{x} < z] < \alpha\} = \inf \{z \mid \mathbb{P}[\tilde{x} \leq z] \geq \alpha\} \\ q_{\tilde{x}}^+(\alpha) &= \inf \{z \mid \mathbb{P}[\tilde{x} \leq z] > \alpha\} = \sup \{z \mid \mathbb{P}[\tilde{x} < z] \leq \alpha\}. \end{aligned}$$

263 Note that when the distribution of  $\tilde{x}$  is absolutely continuous (atomless), then  $q_{\tilde{x}}^+(\alpha) = q_{\tilde{x}}^-(\alpha)$  and  
 264 the quantile is unique. The following example illustrates a simple setting in which the quantile is not  
 265 unique.

266 *Example 1* (Bernoulli random variable). Consider a Bernoulli random variable  $\tilde{e}$  such that  $\tilde{e} = 1$  and  
 267  $\tilde{e} = 0$  with equal (50%) probabilities. Then, any value  $q \in [0, 1]$  is a valid 0.5-quantile because

$$\begin{aligned} q_{\tilde{e}}^-(0.5) &= \inf_{z \in \mathbb{R}} \{z \mid \mathbb{P}[\tilde{e} \leq z] \geq 0.5\} = \inf_{z \in \mathbb{R}} \{z \mid z \geq 0\} = 0 \\ q_{\tilde{e}}^+(0.5) &= \sup_{z \in \mathbb{R}} \{z \mid \mathbb{P}[\tilde{e} \geq z] \geq 0.5\} = \sup_{z \in \mathbb{R}} \{z \mid z \leq 1\} = 1. \end{aligned}$$

268 The objective in Li et al. (2022) is to maximize the quantile operator  $Q_\alpha: \mathbb{X} \rightarrow \mathbb{R}$  defined for a  
 269 reward random variable  $\tilde{x} \in \mathbb{X}$  and a risk level  $\alpha \in [0, 1]$  as

$$Q_\alpha(\tilde{x}) = \inf_{z \in \mathbb{R}} \{z \mid \mathbb{P}[\tilde{x} \leq z] \geq \alpha\}. \quad (16)$$

270 The quantile operator  $Q_\alpha$  and VaR differ in which quantile of the random variable they consider:

$$Q_\alpha(\tilde{x}) = q_{\tilde{x}}^-(\alpha), \quad \text{but} \quad \text{VaR}_\alpha[\tilde{x}] = q_{\tilde{x}}^+(\alpha). \quad (17)$$

271 As a result, the quantile MDP objective in (16) coincides with the VaR value only when the quantile  
 272 is unique, which is not always the case, as shown in Example 1.

273 **Theorem 5.1.** Let  $\tilde{y}: \Omega \rightarrow [N]$  be a random variable distributed as  $\hat{\boldsymbol{p}} = (\hat{p}_i)_{i=1}^N$  with  $\hat{p}_i > 0$ . Then  
 274 for any random variable  $\tilde{x} \in \mathbb{X}$ , we have

$$\text{VaR}_\alpha[\tilde{x}] = \sup_{\zeta \in \Delta_N} \left\{ \min_i \text{VaR}_{\alpha \zeta_i \hat{p}_i^{-1}}[\tilde{x} \mid \tilde{y} = i] \mid \alpha \cdot \zeta \leq \hat{\boldsymbol{p}} \right\}, \quad (18)$$

275 where we interpret the minimum to evaluate to  $\infty$  if all terms are infinite, which only occurs if  $\alpha = 1$ .

276 *Proof.* We first decompose VaR using the definition in (1) as

$$\begin{aligned} \text{VaR}_\alpha[\tilde{x}] &= \sup_{z \in \mathbb{R}} \{z \mid \mathbb{P}[\tilde{x} < z] \leq \alpha\} \stackrel{(a)}{=} \sup_{z \in \mathbb{R}} \left\{ z \mid \sum_{i=1}^N \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \hat{p}_i \leq \alpha \right\} \\ &\stackrel{(b)}{=} \sup_{z \in \mathbb{R}, \zeta \in [0,1]^N} \left\{ z \mid \sum_{i=1}^N \zeta_i \hat{p}_i \leq \alpha, \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \leq \zeta_i, \forall i \in [N] \right\} \\ &\stackrel{(c)}{=} \sup_{z \in \mathbb{R}, \zeta \in [0,1]^N} \left\{ z \mid z \leq \text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i], \forall i \in [N], \sum_{i=1}^N \zeta_i \hat{p}_i \leq \alpha \right\} \\ &\stackrel{(d)}{=} \sup_{\zeta \in [0,1]^N} \left\{ \sup_{z \in \mathbb{R}} \{z \mid z \leq \text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i], \forall i \in [N]\} \mid \sum_{i=1}^N \zeta_i \hat{p}_i \leq \alpha \right\} \\ &\stackrel{(e)}{=} \sup_{\zeta \in [0,1]^N} \left\{ \min_{i \in [N]} \text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i] \mid \sum_{i=1}^N \zeta_i \hat{p}_i \leq \alpha \right\}. \end{aligned}$$

277 We decompose the probability  $\mathbb{P}[\tilde{x} < z]$  in terms of the conditional probabilities  $\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i]$  in  
 278 step (a) and then lower-bound them by an auxiliary variable  $\zeta_i$  in step (b). In step (c), we exploit the  
 279 following equivalence:

$$\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \leq \zeta_i \iff z \leq \text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i]$$

280 The direction  $\Leftarrow$  in the above equivalence follows immediately from the fact that  $\text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i]$  is  
 281 a  $\zeta_i$ -quantile and satisfies (15), i.e.,

$$z \leq \text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i] = q_{\tilde{x}}^+(\zeta_i \mid \tilde{y} = i) \Rightarrow \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \leq \mathbb{P}[\tilde{x} < q_{\tilde{x}}^+(\zeta_i \mid \tilde{y} = i) \mid \tilde{y} = i] \leq \zeta_i.$$

282 The direction  $\Rightarrow$  follows from the definition of VaR (see Eq. 1), which implies that VaR upper-bounds  
 283 any  $z$  that satisfies the left-hand condition:

$$\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \leq \zeta_i \Rightarrow \text{VaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i] = \sup \{z \in \mathbb{R} \mid \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \leq \zeta_i\} \geq z.$$

284 In step (e), we solve for  $z$ . Finally, the form in (18) follows by replacing each  $\zeta_i$  by  $\alpha \zeta_i \hat{p}_i^{-1}$ .  $\square$

285 Focusing on the finite MDP with horizon  $T = 1$ , we can show that the decomposition proposed in  
 286 Theorem 5.1 is amenable to policy optimization. The main difference between the VaR decomposition  
 287 and CVaR is that the former VaR was expressed as a supremum instead of an infimum over quantile  
 288 levels  $\zeta$ . For VaR, changing the order of maximum ( $\pi$ ) and supremum ( $\zeta$ ) does not suffer from a  
 289 potential gap, but changing the order of maximum ( $\pi$ ) and infimum/minimum ( $\zeta$ ) in CVaR does  
 290 suffer from such a gap as shown in Theorem 3.2.

291 The following theorem (proved in Appendix A.4) summarizes the decomposition for VaR.

292 **Theorem 5.2.** *Given any finite MDP with horizon  $T = 1$  and  $\alpha \in [0, 1]$ , we have*

$$\max_{\pi \in \Pi} \text{VaR}_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})}[r(\tilde{s}, \tilde{a}, \tilde{s}')] = \sup_{\zeta \in \Delta_S} \left\{ \min_{s \in \mathcal{S}} \max_{d \in \Delta_A} \text{VaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim d} [r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s] \mid \alpha \cdot \zeta \leq \hat{p} \right\}.$$

293 For completeness, we also present the valid decomposition for the lower quantile MDP (see Appendix  
 294 A.5 for a proof).

295 **Proposition 5.3.** *Given any finite MDP with horizon  $T = 1$  and some  $\alpha \in [0, 1]$ , we have that:*

$$\max_{\pi \in \Pi} Q_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})}(r(\tilde{s}, \tilde{a}, \tilde{s}')) = \sup_{\zeta \in [0, 1]^S} \left\{ \min_{s \in \mathcal{S}: \zeta_s < 1} \max_{d \in \Delta_A} Q_{\zeta_s}^{\tilde{a} \sim d}(r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s) \mid \sum_{s=1}^S \zeta_s \hat{p}_s < \alpha \right\}.$$

296 We note that the difference with the result presented in (Li et al., 2022) resides in the constraint  
 297 imposed on  $\zeta$  that replaces the weak inequality with a strict one. In fact, this strict versus weak  
 298 inequality is the main distinguishing factor between the decompositions for the lower and upper  
 299 quantile.

## 300 6 Conclusion

301 This paper shows that a popular decomposition approach to solving MDPs with CVaR and EVaR  
 302 objectives is suboptimal despite the claims to the contrary. This suboptimality arises from a saddle-  
 303 point gap when *optimizing policy*. We also prove that a similar decomposition approach is optimal for  
 304 policy optimization and evaluation when solving MDPs with the VaR objective. The decomposition  
 305 is optimal because VaR does not involve the same saddle point problem as CVaR and EVaR.

306 Our findings are significant because practitioners who make risk-averse decisions in high-stakes  
 307 scenarios need to have confidence in the correctness of the algorithms they use. Our work raises  
 308 awareness that popular static CVaR and EVaR MDP algorithms are suboptimal, and their analyses  
 309 are inaccurate. We hope the results we present in our paper will increase the scrutiny of dynamic  
 310 programming methods for risk-averse MDPs and motivate research into alternative approaches, such  
 311 as the parametric dynamic programs.

## 312 References

- 313 Mohamadreza Ahmadi, Xiaobin Xiong, and Aaron D Ames. Risk-averse control via CVaR barrier  
314 functions: Application to bipedal robot locomotion. *IEEE Control Systems Letters*, 6:878–883,  
315 2021.
- 316 A. Ahmadi-Javid. Entropic Value-at-Risk: A new coherent risk measure. *Journal of Optimization*  
317 *Theory and Applications*, 155(3):1105–1123, 2012.
- 318 Philippe Artzner, Freddy Delbaen, Jean-Marc Eber, and David Heath. Coherent measures of risk.  
319 *Mathematical Finance*, 9(3):203–228, 1999.
- 320 Basel Committee on Banking Supervision. Minimum capital requirements for market risk. In *Basel*  
321 *III: international regulatory framework for banks*, 2019.
- 322 Nicole Bäuerle and Jonathan Ott. Markov decision processes with average-value-at-risk criteria.  
323 *Mathematical Methods of Operations Research*, 74(3):361–379, 2011.
- 324 Kang Boda, Jerzy A Filar, Yuanlie Lin, and Lieneke Spanjers. Stochastic target hitting time and the  
325 problem of early retirement. *IEEE Transactions on Automatic Control*, 49(3):409–419, 2004.
- 326 Margaret P Chapman, Jonathan Lacotte, Aviv Tamar, Donggun Lee, Kevin M Smith, Victoria  
327 Cheng, Jaime F Fisac, Susmit Jha, Marco Pavone, and Claire J Tomlin. A risk-sensitive finite-  
328 time reachability approach for safety of stochastic dynamic systems. In *2019 American Control*  
329 *Conference (ACC)*, pages 2958–2963. IEEE, 2019.
- 330 Margaret P Chapman, Riccardo Bonalli, Kevin M Smith, Insoon Yang, Marco Pavone, and Claire J  
331 Tomlin. Risk-sensitive safety analysis using conditional value-at-risk. *IEEE Transactions on*  
332 *Automatic Control*, 67(12):6521–6536, 2022.
- 333 Arnav Choudhry, Brady Moon, Jay Patrikar, Constantine Samaras, and Sebastian Scherer. CVaR-  
334 based flight energy risk assessment for multirotor uavs using a deep energy model. In *2021 IEEE*  
335 *International Conference on Robotics and Automation (ICRA)*, pages 262–268. IEEE, 2021.
- 336 Yinlam Chow and Mohammad Ghavamzadeh. Algorithms for CVaR optimization in MDPs. In  
337 *Neural Information Processing Systems (NIPS)*, 2014.
- 338 Yinlam Chow, Aviv Tamar, Shie Mannor, and Marco Pavone. Risk-sensitive and robust decision-  
339 making: A CVaR optimization approach. In *Neural Information Processing Systems (NIPS)*,  
340 2015.
- 341 Yinlam Chow, Mohammad Ghavamzadeh, Lucas Janson, and Marco Pavone. Risk-constrained  
342 reinforcement learning with percentile risk criteria. *Journal of Machine Learning Research*, 18:  
343 1–51, 2018.
- 344 Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. Wiley-Interscience, 2<sup>nd</sup>  
345 edition, 2006.
- 346 Erick Delage, Daniel Kuhn, and Wolfram Wiesemann. “Dice”-sion-making under uncertainty: When  
347 can a random decision reduce risk? *Management Science*, 65(7):3282–3301, 2019.
- 348 Rui Ding and Eugene Feinberg. CVaR optimization for MDPs: Existence and computation of optimal  
349 policies. *ACM SIGMETRICS Performance Evaluation Review*, 50(2):39–41, 2022.
- 350 Jerzy A. Filar, Dmitry Krass, and Keith W. Ross. Percentile Performance Criteria For Limiting  
351 Average Markov Decision Processes. *IEEE Transactions on Automatic Control*, 40(1):2–10, 1995.
- 352 Hans Follmer and Alexander Schied. *Stochastic Finance: Introduction in Discrete Time*. De Gruyter  
353 Graduate, fourth edition, 2016.
- 354 Astghik Hakobyan and Insoon Yang. Wasserstein distributionally robust motion control for collision  
355 avoidance using conditional value-at-risk. *IEEE Transactions on Robotics*, 38(2):939–957, 2021.
- 356 Jia Lin Hau, Marek Petrik, and Mohammad Ghavamzadeh. Entropic risk optimization in discounted  
357 MDPs. In *Artificial Intelligence and Statistics (AISTATS)*, 2023.

- 358 I Ge Jin, Bastian Schürmann, Richard M Murray, and Matthias Althoff. Risk-aware motion planning  
359 for automated vehicle among human-driven cars. In *2019 American Control Conference (ACC)*,  
360 pages 3987–3993. IEEE, 2019.
- 361 Ramtin Keramati, Christoph Dann, Alex Tamkin, and Emma Brunskill. Being optimistic to be  
362 conservative: Quickly learning a CVaR policy. In *Proceedings of the AAAI conference on artificial*  
363 *intelligence*, volume 34, pages 4436–4443, 2020.
- 364 Ümit Emre Köse. *Optimal timing of living-donor liver transplantation under risk-aversion*. PhD  
365 thesis, Bilkent Üniversitesi (Turkey), 2016.
- 366 Xiaocheng Li, Huaiyang Zhong, and Margaret L. Brandeau. Quantile Markov decision processes.  
367 *Operations Research*, 70(3):1428–1447, 2022.
- 368 Yuanlie Lin, Congbin Wu, and Boda Kang. Optimal models with maximizing probability of first  
369 achieving target value in the preceding stages. *Science in China Series A: Mathematics*, 46:  
370 396–414, 2003.
- 371 Seungki Min, Ciamac C Moallemi, and Costis Maglaras. Risk-sensitive optimal execution via a  
372 conditional value-at-risk objective. *arXiv preprint arXiv:2201.11962*, 2022.
- 373 Xinyi Ni and Lifeng Lai. Risk-sensitive reinforcement learning via Entropic-VaR optimization. In  
374 *2022 56<sup>th</sup> Asilomar Conference on Signals, Systems, and Computers*, pages 953–959. IEEE, 2022.
- 375 Georg Ch Pflug and Alois Pichler. Time-consistent decisions and temporal decomposition of coherent  
376 risk functionals. *Mathematics of Operations Research*, 41(2):682–699, 2016.
- 377 L. A Prashanth and Michael Fu. Risk-sensitive reinforcement learning via policy gradient search.  
378 *Foundations and Trends in Machine Learning*, 15(5):537–693, 2022.
- 379 Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*.  
380 Wiley-Interscience, 2005.
- 381 Marc Rigter, Bruno Lacerda, and Nick Hawes. Risk-averse Bayes-adaptive reinforcement learning.  
382 *Advances in Neural Information Processing Systems*, 34:1142–1154, 2021.
- 383 A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and*  
384 *Theory*. SIAM, 2014.
- 385 Alexander Shapiro. Interchangeability principle and dynamic equations in risk averse stochastic  
386 programming. *Operations Research Letters*, 45(4):377–381, 2017.
- 387 Vishnu D Sharma, Maymoonah Toubeh, Lifeng Zhou, and Pratap Tokekar. Risk-aware planning and  
388 assignment for ground vehicles using uncertain perception from aerial vehicles. In *2020 IEEE/RSJ*  
389 *International Conference on Intelligent Robots and Systems (IROS)*, pages 11763–11769. IEEE,  
390 2020.
- 391 Silvestr Stanko and Karel Macek. Risk-averse distributional reinforcement learning: A CVaR  
392 optimization approach. In *International Joint Conference on Computational Intelligence*, pages  
393 412–423, 2019.
- 394 Congbin Wu and Yuanlie Lin. Minimizing risk models in Markov decision processes with policies  
395 depending on target values. *Journal of mathematical analysis and applications*, 231(1):47–67,  
396 1999.
- 397 Huan Xu and Shie Mannor. Probabilistic goal Markov decision processes. In *International Joint*  
398 *Conference on Artificial Intelligence*, 2011.
- 399 Huaiyang Zhong. *Decision Making for Disease Treatment: Operations Research and Data Analytic*  
400 *Modeling*. Stanford University, 2020.

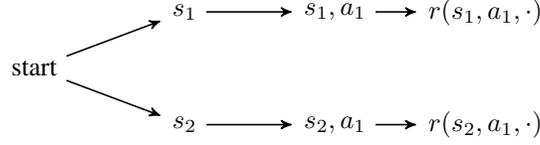


Figure 3: Rewards of the MDP  $M_E$  used in the proof of Theorem 4.1. The dot indicates that the rewards are independent of the next state.

## 401 A Proofs

### 402 A.1 Proof of Proposition 3.1

403 Suppose that  $\alpha > 0$ ; the decomposition for  $\alpha = 0$  holds readily because  $\text{CVaR}_0[\tilde{x}] = \text{ess inf}[\tilde{x}]$ .

404 To streamline the notation, Define a random variable  $\tilde{x} = r(\tilde{s}, \tilde{a}, \tilde{s}')$  over a random space  $\Omega =$   
 405  $\mathcal{S} \times \mathcal{A} \times \mathcal{S}$  with a probability distribution  $\mathbf{q} \in \Delta_m$  such that  $q_{s,a,s'} = \hat{p}_s \cdot \pi(s,a) \cdot p(s,a,s')$ . The  
 406 value  $\mathbf{x}$  is the vector representation of the random variable  $\tilde{x}$  and  $\boldsymbol{\xi}_s = \boldsymbol{\xi}_{s,\cdot,\cdot} \in \mathbb{R}^{\mathcal{S} \cdot \mathcal{A}}$  for  $\boldsymbol{\xi} \in \mathbb{R}^m$  is a  
 407 vector that corresponds to the subset of the elements of  $\Omega$  in which the first element is some  $s \in \mathcal{S}$ .  
 408 The vectors  $\mathbf{x}_s = \mathbf{x}_{s,\cdot,\cdot} \in \mathbb{R}^{\mathcal{S} \cdot \mathcal{A}}$  and  $\mathbf{q}_s = \mathbf{q}_{s,\cdot,\cdot} \in \mathbb{R}^{\mathcal{S} \cdot \mathcal{A}}$  are defined analogously to  $\boldsymbol{\xi}_s$ .

409 Starting with the CVaR definition in (2) and introducing an auxiliary variable  $\zeta$  we get that

$$\begin{aligned} \text{CVaR}_\alpha[\tilde{x}] &= \min_{\boldsymbol{\xi} \in \Delta_m} \{ \mathbf{x}^\top \boldsymbol{\xi} \mid \alpha \boldsymbol{\xi} \leq \mathbf{q} \} = \min_{\boldsymbol{\xi} \in \Delta_m, \zeta \in \mathbb{R}^{\mathcal{S}}} \{ \mathbf{x}^\top \boldsymbol{\xi} \mid \alpha \boldsymbol{\xi} \leq \mathbf{q}, \zeta_s = \mathbf{1}^\top \boldsymbol{\xi}_s, \forall s \in \mathcal{S} \} \\ &= \min_{\boldsymbol{\xi} \in \Delta_m, \zeta \in \Delta_{\mathcal{S}}} \{ \mathbf{x}^\top \boldsymbol{\xi} \mid \alpha \boldsymbol{\xi} \leq \mathbf{q}, \zeta_s = \mathbf{1}^\top \boldsymbol{\xi}_s, \alpha \zeta_s \leq \hat{p}_s, \forall s \in \mathcal{S} \} \\ &= \min_{\boldsymbol{\xi} \in \mathbb{R}_+^\Omega, \zeta \in \mathcal{Z}_C} \{ \mathbf{x}^\top \boldsymbol{\xi} \mid \alpha \boldsymbol{\xi} \leq \mathbf{q}, \zeta_s = \mathbf{1}^\top \boldsymbol{\xi}_s, \forall s \in \mathcal{S} \}. \end{aligned}$$

410 In the derivation above, we replaced the infimum by a minimum because  $\Omega$  is finite, introduced a  
 411 new variable  $\zeta$ , derived implied constraints on  $\zeta$ , and then dropped superfluous constraints on  $\boldsymbol{\xi}$ .  
 412 Continuing with the derivation above and noticing that the constraints on each  $\boldsymbol{\xi}_s$  are independent  
 413 given  $\zeta$ , we get that

$$\begin{aligned} \text{CVaR}_\alpha[\tilde{x}] &= \min_{\boldsymbol{\xi} \in \mathbb{R}_+^\Omega, \zeta \in \mathcal{Z}_C} \left\{ \sum_{s \in \mathcal{S}} \mathbf{x}_s^\top \boldsymbol{\xi}_s \mid \alpha \boldsymbol{\xi} \leq \mathbf{q}, \zeta_s = \mathbf{1}^\top \boldsymbol{\xi}_s, \forall s \in \mathcal{S} \right\} \\ &\stackrel{(a)}{=} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \inf_{\boldsymbol{\xi}_s \in \mathbb{R}_+^{\Omega_s}} \{ \mathbf{x}_s^\top \boldsymbol{\xi}_s \mid \alpha \boldsymbol{\xi}_s \leq \mathbf{q}_s, \zeta_s = \mathbf{1}^\top \boldsymbol{\xi}_s \} \\ &\stackrel{(b)}{=} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \cdot \min_{\boldsymbol{\chi} \in \Delta_{\mathcal{S} \cdot \mathcal{A}}} \{ \mathbf{x}_s^\top \boldsymbol{\chi} \mid \alpha \hat{p}_s^{-1} \zeta_s \boldsymbol{\chi}_{a,s'} \leq \hat{p}_s^{-1} q_{s,a,s'}, \forall a \in \mathcal{A}, s' \in \mathcal{S} \} \\ &\stackrel{(c)}{=} \min_{\zeta \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \zeta_s \cdot \text{CVaR}_{\alpha \zeta_s \hat{p}_s^{-1}}[\tilde{x} \mid \tilde{s} = s]. \end{aligned}$$

414 The step (a) follows from the interchangeability principle (Shapiro et al., 2014, theorem 7.92), and the  
 415 step (b) follows by substituting  $\boldsymbol{\xi}_{s,a,s'} = \zeta_s \boldsymbol{\chi}_{a,s'}$  taking care when  $\zeta_s = 0$  and multiplying both sides  
 416 of the inequality by  $\hat{p}_s^{-1} > 0$ . Finally, in step (c), the random variable  $\tilde{x} = r(\tilde{s}, \tilde{a}, \tilde{s}')$  conditional  
 417 on  $\tilde{s} = s$  is distributed according to  $q_{s,a,s'} \hat{p}_s^{-1}$  and the equality follows from the definition of CVaR  
 418 in (2).  $\square$

### 419 A.2 Proof of Theorem 4.1

420 Consider an MDP  $M_E$  depicted in Figure 3 with  $\mathcal{S} = \{s_1, s_2\}$  and  $\mathcal{A} = \{a_1\}$  and a reward function  
 421  $r(s_1, a_1, \cdot) = 1$  and  $r(s_2, a_1, \cdot) = 0$ . We abbreviate the rewards to  $r(s_1)$  and  $r(s_2)$  because they only  
 422 depend on the originating state. The initial distribution is  $\hat{p}_{s_1} = \hat{p}_{s_2} = 0.5$ . We finally let  $\alpha = 0.75$ .

423 Because  $\mathcal{Z}_E \subseteq \mathcal{Z}_C$ , the right-hand side of (13) can be lower-bounded by CVaR as

$$\begin{aligned} \min_{\boldsymbol{\xi} \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}}[r(s, a_1, \tilde{s}')] &= \min_{\boldsymbol{\xi} \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s r(s) \\ &\geq \min_{\boldsymbol{\xi} \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \xi_s r(s) = \text{CVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] . \end{aligned} \tag{19}$$

424 The first equality holds from the positive homogeneity and cash invariance properties of EVaR, and  
 425 the last equality follows from the dual representation of CVaR (Follmer and Schied, 2016).

426 Because  $\text{EVaR}_\alpha[\tilde{x}] \leq \text{CVaR}_\alpha[\tilde{x}]$  for each  $\alpha \in [0, 1]$  and  $\tilde{x} \in \mathbb{X}$  (see (Ahmadi-Javid, 2012,  
 427 Proposition 3.2)), we can further lower-bound (19) as

$$\text{EVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] \leq \text{CVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] \leq \min_{\xi \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}}[r(s, a_1, \tilde{s}')] . \quad (20)$$

428 Therefore, (13) holds with an inequality.

429 To prove by contradiction that the inequality in (13) is strict, suppose that

$$\text{EVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] = \min_{\xi \in \mathcal{Z}_E} \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}}[r(s, a_1, \tilde{s}')] . \quad (21)$$

430 Equalities (21) and (20) imply that  $\text{EVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] = \text{CVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')]$  which is false in  
 431 general (Ahmadi-Javid, 2012).

432 We now show that EVaR does not equal CVaR even for the categorical distribution of  $\tilde{s}$ . The CVaR of  
 433 the return in  $M_E$  reduces from (2) to

$$\text{CVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] = \min_{\xi \in \mathcal{Z}_C} \sum_{s \in \mathcal{S}} \xi_s r(s) = \max \left\{ 0, \frac{\hat{p}_{s_1} + \alpha - 1}{\alpha} \right\} . \quad (22)$$

434 Since  $1 - \alpha = 0.25 < 0.5 = \hat{p}_{s_1}$ , then the optimal  $\xi^*$  in (22) is

$$\xi^* = \left( \frac{\hat{p}_{s_1} + \alpha - 1}{\frac{1 - \alpha}{\hat{p}_{s_1}}} \right) .$$

435 Since  $\text{KL}(\xi^* \|\hat{p}) < -\log \alpha$ , we have that  $\xi^*$  is in the relative interior of the EVaR feasible region  
 436 in (3), and, therefore, there exists an  $\epsilon > 0$  such that

$$\text{EVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] = \text{CVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] - \epsilon < \text{CVaR}_\alpha[r(\tilde{s}, a_1, \tilde{s}')] ,$$

437 which proves the desired inequality. □

### 438 A.3 Proof of Theorem 4.2

439 We start by proposing a new decomposition for EVaR.

440 **Proposition A.1.** *Given a random variable  $\tilde{x} \in \mathbb{X}$  and a discrete variable  $\tilde{y}: \Omega \rightarrow \mathcal{N} = \{1, \dots, N\}$ ,  
 441 with probabilities denoted as  $\{\hat{p}_i\}_{i=1}^N$ , for any  $\alpha \in (0, 1]$  we have that*

$$\text{EVaR}_\alpha[\tilde{x}] = \inf_{\zeta \in (0, 1]^N} \min_{\xi \in \mathcal{Z}'_E(\zeta)} \sum_i \xi_i \text{EVaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i] ,$$

442 where

$$\mathcal{Z}'_E(\zeta) = \left\{ \xi \in \Delta_N \mid \xi \ll \hat{p}, \sum_{i=1}^N \xi_i (\log(\xi_i / \hat{p}_i) - \log(\zeta_i)) \leq -\log \alpha \right\} .$$

443 *Proof.* Let  $\mathbf{q}$  denote the joint probability distribution of  $\tilde{x}$  and  $\tilde{y}$ . The proof exploits the chain rule of  
 444 relative entropy (e.g., Cover and Thomas (2006, theorem 2.5.3)), which states that for any probability  
 445 distributions  $\boldsymbol{\eta}, \mathbf{q} \in \Delta_\Omega$  with  $\boldsymbol{\eta} \ll \mathbf{q}$

$$\text{KL}(\boldsymbol{\eta} \|\mathbf{q}) = \text{KL}(\boldsymbol{\eta}(\tilde{y}) \|\mathbf{q}(\tilde{y})) + \text{KL}(\boldsymbol{\eta}(\tilde{x}|\tilde{y}) \|\mathbf{q}(\tilde{x}|\tilde{y})), \quad (23)$$

446 where the conditional relative entropy is defined as

$$\text{KL}(\boldsymbol{\eta}(\tilde{x}|\tilde{y}) \|\mathbf{q}(\tilde{x}|\tilde{y})) = \mathbb{E}^\boldsymbol{\eta} \left[ \log \frac{\boldsymbol{\eta}(\tilde{x}|\tilde{y})}{\mathbf{q}(\tilde{x}|\tilde{y})} \right] .$$

447 with  $\mathbb{E}^\eta[f(\tilde{x}, \tilde{y})]$  as a shorthand notation to indicate that  $(\tilde{x}, \tilde{y}) \sim \eta$ . We can now decompose EVaR  
 448 from its definition in (3) as

$$\begin{aligned}
 \text{EVaR}_\alpha[\tilde{x}] &= \inf_{\eta \in \Delta_{m,N} : \eta \ll \mathbf{q}} \left\{ \mathbb{E}^\eta[\tilde{x}] \mid \text{KL}(\eta \mid \mathbf{q}) \leq -\log \alpha \right\} \\
 &\stackrel{(a)}{=} \inf_{\eta \in \Delta_{m,N} : \eta \ll \mathbf{q}} \left\{ \mathbb{E}^\eta[\tilde{x}] \mid \text{KL}(\eta(\tilde{y}) \parallel \mathbf{q}(\tilde{y})) + \text{KL}(\eta(\tilde{x} \mid \tilde{y}) \parallel \mathbf{q}(\tilde{x} \mid \tilde{y})) \leq -\log \alpha \right\} \\
 &= \inf_{\eta \in \Delta_{m,N} : \eta \ll \mathbf{q}} \left\{ \mathbb{E}^\eta[\tilde{x}] \mid \text{KL}(\eta(\tilde{y}) \parallel \mathbf{q}(\tilde{y})) + \mathbb{E}^\eta \left[ \mathbb{E}^\eta \left[ \log \frac{\eta(\tilde{x} \mid \tilde{y})}{\mathbf{q}(\tilde{x} \mid \tilde{y})} \mid \tilde{y} \right] \leq -\log \alpha \right\} \\
 &\stackrel{(b)}{=} \inf_{\eta \in \Delta_{m,N}, \zeta \in (0,1]^N : \eta \ll \mathbf{q}} \left\{ \mathbb{E}^\eta[\mathbb{E}^\eta[\tilde{x} \mid \tilde{y}]] \mid \frac{\text{KL}(\eta(\tilde{y}) \parallel \mathbf{q}(\tilde{y})) + \mathbb{E}^\eta[-\log(\zeta_{\tilde{y}})] \leq -\log \alpha}{\mathbb{P}_\eta[\mathbb{E}^\eta[\log(\eta(\tilde{x} \mid \tilde{y})/\mathbf{q}(\tilde{x} \mid \tilde{y})) \mid \tilde{y}] \leq -\log(\zeta_{\tilde{y}})] = 1} \right\} \\
 &\stackrel{(c)}{=} \inf_{\xi \in \Delta_N, \zeta \in (0,1]^N : \xi \ll \hat{\mathbf{p}}} \left\{ \mathbb{E}^\xi[\text{EVaR}_{\zeta_{\tilde{y}}}[\tilde{x} \mid \tilde{y}]] \mid \text{KL}(\xi \parallel \hat{\mathbf{p}}) + \mathbb{E}^\xi[-\log(\zeta_{\tilde{y}})] \leq -\log \alpha \right\} \\
 &= \inf_{\xi \in \Delta_N, \zeta \in (0,1]^N : \xi \ll \hat{\mathbf{p}}} \left\{ \sum_i \xi_i \text{EVaR}_{\zeta_i}[\tilde{x} \mid \tilde{y} = i] \mid \sum_{i=1}^N \xi_i \log(\xi_i/\hat{p}_i) - \sum_{i=1}^N \xi_i \log(\zeta_i) \leq -\log \alpha \right\}.
 \end{aligned}$$

449 Here, we decompose the relative entropy of  $\eta$  and  $\mathbf{q}$  using (23) in step (a) and then use the tower  
 450 property of the expectation operator in the next step. In step (b), we introduce a variable  $\zeta_i$  for each  
 451 realization of  $\tilde{y} = i$  with  $i \in \mathcal{N}$  to decouple the influence of  $\eta(\tilde{x} \mid \tilde{y})$ , under each  $\tilde{y}$ , in the inequality  
 452 constraint. Finally, we replace the conditional EVaR definition by solving for  $\eta(\tilde{x} \mid \tilde{y})$  for a given  $\zeta$   
 453 in step (c), and representing  $\eta(\tilde{y})$  using  $\xi$ .  $\square$

454 The first part of our theorem follows directly from Proposition A.1. Suppose that  $\alpha > 0$ ; the result  
 455 follows for  $\alpha = 0$  because  $\text{EVaR}_0[\cdot]$  reduces to  $\text{ess inf}$ . Then, the second part of the corollary holds  
 456 as

$$\begin{aligned}
 \text{EVaR}_\alpha[r(\tilde{s}, \tilde{a}, \tilde{s}')] &= \inf_{\zeta \in (0,1]^N, \xi \in \Delta_N} \left\{ \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\zeta_s}[r(s, \tilde{a}, \tilde{s}')] \mid \sum_{s \in \mathcal{S}} \xi_s \log \frac{\xi_s}{\zeta_s \hat{p}_s} \leq -\log \alpha \right\} \\
 &\leq \inf_{\zeta \in (0,1]^N, \xi \in \Delta_N} \left\{ \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\zeta_s}[r(s, \tilde{a}, \tilde{s}')] \mid \sum_{s \in \mathcal{S}} \xi_s \log \frac{\xi_s}{\zeta_s \hat{p}_s} \leq -\log \alpha, \xi \leq \alpha^{-1} \hat{\mathbf{p}} \right\} \\
 &\leq \inf_{\xi \in \Delta_N} \left\{ \sum_{s \in \mathcal{S}} \xi_s \text{EVaR}_{\alpha \xi_s \hat{p}_s^{-1}}[r(s, \tilde{a}, \tilde{s}')] \mid \xi \leq \alpha^{-1} \hat{\mathbf{p}} \right\}.
 \end{aligned}$$

457 The first inequality follows from adding a constraint on the pairs on the  $\xi$  considered by the infimum.  
 458 The second inequality follows by fixing  $\zeta_s = \hat{\zeta}_s$  with  $\hat{\zeta}_s = \alpha \xi_s \hat{p}_s^{-1}$  for each  $s \in \mathcal{S}$ . This is an upper  
 459 bound because  $\hat{\zeta}_s$  is feasible in the infimum:

$$\sum_{s \in \mathcal{S}} \xi_s \log \frac{\xi_s}{\hat{\zeta}_s \hat{p}_s} = -\log \alpha \leq -\log \alpha.$$

460 The value  $\hat{\zeta}_s$  is well-defined since  $\hat{p}_s > 0$  and the constraint  $\xi \leq \alpha^{-1} \hat{\mathbf{p}}$  ensures that  $\hat{\zeta}_s \leq 1$ . Also, we  
 461 can relax the constraint  $\zeta_s > 0 \Rightarrow \xi_s > 0$  to  $\xi_s \geq 0$  because  $\text{EVaR}_0[\tilde{x}] = \lim_{\alpha \rightarrow 0} \text{EVaR}_\alpha[\tilde{x}]$ , and,  
 462 therefore, the infimum is not affected. Finally, the inequality in the corollary follows immediately by  
 463 further upper bounding the decomposition above by adding a constraint.  $\square$

#### 464 A.4 Proof of Theorem 5.2

465 The equality develops from Theorem 5.1 as

$$\begin{aligned}
 \max_{\pi \in \Pi} \text{VaR}_\alpha^{\tilde{a} \sim \pi(\tilde{s})}[r(\tilde{s}, \tilde{a}, \tilde{s}')] &= \max_{\pi \in \Pi} \sup_{\zeta \in \Delta_S : \alpha \cdot \zeta \leq \hat{\mathbf{p}}} \min_{s \in \mathcal{S}} \left( \text{VaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(s)}[r(s, \tilde{a}, \tilde{s}')] \mid \tilde{s} = s \right) \\
 &= \sup_{\zeta \in \Delta_S : \alpha \cdot \zeta \leq \hat{\mathbf{p}}} \max_{\pi \in \Pi} \min_{s \in \mathcal{S}} \left( \text{VaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim \pi(s)}[r(s, \tilde{a}, \tilde{s}')] \mid \tilde{s} = s \right) \\
 &= \sup_{\zeta \in \Delta_S : \alpha \cdot \zeta \leq \hat{\mathbf{p}}} \min_{s \in \mathcal{S}} \left( \max_{d \in \Delta_A} \text{VaR}_{\alpha \zeta_s \hat{p}_s^{-1}}^{\tilde{a} \sim d} [r(s, \tilde{a}, \tilde{s}')] \mid \tilde{s} = s \right),
 \end{aligned}$$

466 where we first change the order of maximum and supremum, followed by changing the order of  
 467  $\max_{\pi} \min_s$  with  $\min_s \max_{\pi}$ . The latter is a direct consequence of the interchangeability property of  
 468 the maximum operation (Shapiro, 2017, Proposition 2.2).  $\square$

### 469 A.5 Proof of Proposition 5.3

470 The proof mainly relies on correcting the decomposition of lower quantile proposed in (Li et al.,  
 471 2022).

472 **Proposition A.2.** *Given an  $\tilde{x} \in \mathbb{X}$ , suppose that a random variable  $\tilde{y}: \Omega \rightarrow \mathcal{N} = \{1, \dots, N\}$  is  
 473 distributed as  $\hat{\mathbf{p}} = (\hat{p}_i)_{i=1}^N$  with  $\hat{p}_i > 0$ . Then:*

$$Q_{\alpha}(\tilde{x}) = \sup_{\zeta \in [0,1]^N} \left\{ \min_{i \in \mathcal{N}: \zeta_i < 1} Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i) \mid \sum_{i=1}^N \zeta_i \hat{p}_i < \alpha \right\}, \quad (24)$$

474 where we interpret the supremum to be minus infinity if its feasible set is empty, which only occurs if  
 475  $\alpha = 0$ .

476 *Proof.* First, we decompose lower quantile using its definition as

$$\begin{aligned} Q_{\alpha}(\tilde{x}) &= \sup \{z \mid \mathbb{P}[\tilde{x} < z] < \alpha\} \stackrel{(a)}{=} \sup_{z \in \mathbb{R}} \left\{ z \mid \sum_{i=1}^N \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] \hat{p}_i < \alpha \right\} \\ &\stackrel{(b)}{=} \sup_{z \in \mathbb{R}, \zeta \in [0,1]^N} \left\{ z \mid \sum_{i=1}^N \zeta_i \hat{p}_i < \alpha, \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] < \zeta_i, \forall i \in \mathcal{N} : \zeta_i < 1 \right\} \\ &\stackrel{(c)}{=} \sup_{z \in \mathbb{R}, \zeta \in [0,1]^N} \left\{ z \mid z < Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i), \forall i \in \mathcal{N} : \zeta_i < 1, \sum_{i=1}^N \zeta_i \hat{p}_i < \alpha \right\} \\ &\stackrel{(d)}{=} \sup_{\zeta \in [0,1]^N} \left\{ \sup_{z \in \mathbb{R}} \{z \mid z < Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i), \forall i \in \mathcal{N} : \zeta_i < 1\} \mid \sum_{i=1}^N \zeta_i \hat{p}_i < \alpha \right\} \\ &\stackrel{(e)}{=} \sup_{\zeta \in [0,1]^N} \left\{ \min_{i \in \mathcal{N}: \zeta_i < 1} Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i) \mid \sum_{i=1}^N \zeta_i \hat{p}_i < \alpha \right\}. \end{aligned}$$

477 We decompose the probability  $\mathbb{P}[\tilde{x} < z]$  in terms of the conditional probabilities  $\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i]$  in  
 478 step (a) and then lower-bound them by an auxiliary variable  $\zeta_i$  in step (b). In step (c), we exploits the  
 479 following equivalence:

$$\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] < \zeta_i \Leftrightarrow z < Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i)$$

480 The direction  $\Leftarrow$  in the equivalence follows from the definition of  $Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i)$ :

$$z < Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i) = \inf \{z \mid \mathbb{P}[\tilde{x} \leq z \mid \tilde{y} = i] \geq \zeta_i\} \Rightarrow \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] < \zeta_i.$$

481 The direction  $\Rightarrow$  follows from the definition of VaR (see equation (1)), which implies that VaR  
 482 upper-bounds any  $z$  that satisfies the left-hand condition:

$$\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] < \zeta_i \Rightarrow Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i) = \sup \{z \in \mathbb{R} \mid \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] < \zeta_i\} \geq z,$$

483 yet  $Q_{\zeta_i}(\tilde{x} \mid \tilde{y} = i) \neq z$  otherwise since  $\mathbb{P}[\tilde{x} < z \mid \tilde{y} = i]$  is right continuous, there must exist some  
 484  $\epsilon > 0$  for which  $\mathbb{P}[\tilde{x} < z + \epsilon \mid \tilde{y} = i] < \zeta_i$  hence:

$$z = \sup \{z \in \mathbb{R} \mid \mathbb{P}[\tilde{x} < z \mid \tilde{y} = i] < \zeta_i\} \geq z + \epsilon > z,$$

485 which leads to a contradiction. In step (e), we solve for  $z$ . Finally, we obtain the form in (24).  $\square$

486 The decomposition proposed in Proposition A.2 can now be used, exactly as was done for the case of  
 487 VaR, to obtain a decomposition for the risk averse MDP:

$$\begin{aligned}
 \max_{\pi \in \Pi} Q_{\alpha}^{\tilde{a} \sim \pi(\tilde{s})}[r(\tilde{s}, \tilde{a}, \tilde{s}')] &= \max_{\pi \in \Pi} \sup_{\zeta \in [0,1]^S} \left\{ \min_{s \in \mathcal{S}: \zeta_s < 1} Q_{\zeta_s}^{\tilde{a} \sim \pi(s)}(r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s) \mid \sum_{s=1}^S \zeta_s \hat{p}_s < \alpha \right\} \\
 &= \sup_{\zeta \in [0,1]^S} \max_{\pi \in \Pi} \left\{ \min_{s \in \mathcal{S}: \zeta_s < 1} Q_{\zeta_s}^{\tilde{a} \sim \pi(s)}(r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s) \mid \sum_{s=1}^S \zeta_s \hat{p}_s < \alpha \right\} \\
 &= \sup_{\zeta \in [0,1]^S} \left\{ \min_{s \in \mathcal{S}: \zeta_s < 1} \left( \max_{d \in \Delta_A} Q_{\zeta_s}^{\tilde{a} \sim d}(r(s, \tilde{a}, \tilde{s}') \mid \tilde{s} = s) \mid \sum_{s=1}^S \zeta_s \hat{p}_s < \alpha \right) \right\}.
 \end{aligned}$$