Robust ϕ **-Divergence MDPs**

Anonymous Author(s) Affiliation Address email

Abstract

1	In recent years, robust Markov decision processes (MDPs) have emerged as a
2	prominent modeling framework for dynamic decision problems affected by un-
3	certainty. In contrast to classical MDPs, which only account for stochasticity by
4	modeling the dynamics through a stochastic process with a known transition kernel,
5	robust MDPs additionally account for <i>ambiguity</i> by optimizing in view of the most
6	adverse transition kernel from a prescribed ambiguity set. In this paper, we develop
7	a novel solution framework for robust MDPs with s-rectangular ambiguity sets that
8	decomposes the problem into a sequence of robust Bellman updates and simplex
9	projections. Exploiting the rich structure present in the simplex projections corre-
10	sponding to ϕ -divergence ambiguity sets, we show that the associated s-rectangular
11	robust MDPs can be solved substantially faster than with state-of-the-art commer-
12	cial solvers as well as a recent first-order solution scheme, thus rendering them
13	attractive alternatives to classical MDPs in practical applications.

14 **1** Introduction

Markov decision processes (MDPs) are a flexible and popular framework for dynamic decisionmaking problems and reinforcement learning [40, 50]. A practical limitation of the standard MDP model is that it assumes the model parameters, such as transition probabilities and rewards, to be known exactly. In reinforcement learning and other applications, these parameters must be estimated from sampled data, which introduces estimation errors. Optimal MDP solutions, referred to as policies, are well known to be sensitive to errors and may fail catastrophically when deployed [26, 58].

Robust MDPs (RMDPs) mitigate the sensitivity of MDPs to estimation errors by computing a policy 21 that is optimal for the worst plausible realization of the transition probabilities. This set of plausible 22 transition probabilities is known as the *ambiguity set*. Most prior work considers ambiguity sets that 23 are rectangular. In this work, we focus on s-rectangular ambiguity sets, which assume that the worst 24 transition probabilities are chosen independently in each state [26, 58]. While several other models of 25 rectangularity have been studied [9, 14, 22, 29], s-rectangular ambiguity sets are popular due to their 26 generality and the existence of polynomial-time algorithms based on dynamic programming concepts. 27 It has been shown that s-rectangular sets can provide policies that are less conservative compared to 28 (s, a)-rectangular sets [58], both in-sample and out-of-sample. However, even those algorithms may 29 be too slow in practice. Solving RMDPs requires the solution of a convex optimization problem in 30 every step of value or policy iteration, which can become prohibitively slow even in moderatly sized 31 problems with 100s of states [5, 9, 15, 20]. 32

³³ Motivated by the difficulty of solving RMDPs, several fast algorithms have been proposed for *s*-³⁴ rectangular RMDPs [5, 9, 15, 20]. The preponderance of the earlier work has focused on ambiguity ³⁵ sets defined in terms of L_1 - and L_{∞} -norms. These ambiguity sets are polyhedral, and they can be ³⁶ analyzed using linear programming techniques which offer fruitful avenues to exploit the structure ³⁷ inherent to those sets. However, recent statistical studies point to the superior solution quality offered ³⁸ by nonlinear ambiguity sets defined in terms of the Kullback-Leibler (KL) divergence, the L_2 -norm

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

and other metrics [18]. Linear optimization solvers are not applicable to RMDPs with *s*-rectangular ambiguity sets defined in terms of non-polyhedral ambiguity sets, as the corresponding optimization

41 problems are in general convex conic programs (e.g., exponential cone program in the case of KL

42 divergence); thus, they are currently solved using first-order methods [15] or general convex conic

43 solvers such as MOSEK [3], which tend to be complex, closed-source and slow.

As our main contribution, we propose a new suite of fast algorithms for solving RMDPs with ϕ -44 divergence constrained s-rectangular ambiguity sets. ϕ -divergences, also known as f-divergences, 45 constitute a generalization of the KL divergence that encompasses the Burg entropy as well as 46 the L_1 - and weighted L_2 -norms as special cases [4, 6]. Moreover, ϕ -divergence ambiguity sets 47 benefit from rigorous statistical performance guarantees, and they are optimal among all (known 48 and unknown) data-driven optimization paradigms for certain types of worst-case out-of-sample 49 50 performance guarantees [36]. The radii of ϕ -divergence ambiguity sets can be selected either via crossvalidation or via statistical bounds [27, 33, 57]. Robust MDPs with ϕ -divergence sets is challenging 51 52 and unexplored for both (s, a)-rectangular and s-rectangular ambiguity sets. Solving ϕ -divergence RMDPs using value iteration requires the solution of seemingly unstructured min-max problems. Our 53 main insight is that these min-max problems can be reduced to a small number of highly structured 54 projection problems onto a probability simplex. We use this insight to develop tailored solution 55 schemes for the projection problems corresponding to several popular ϕ -divergence ambiguity sets, 56 which in turn give rise to efficient solution methods for the respective RMDPs. Ignoring tolerances, 57 our algorithms achieve an overall $\mathcal{O}(S^2 \cdot A \log A)$ or $\mathcal{O}(S^2 \log S \cdot A)$ time complexity to compute 58 the robust Bellman operator, where S and A denote the numbers of states and actions, respectively. 59 Since the evaluation of a non-robust Bellman operator requires a runtime of $\mathcal{O}(S^2 \cdot \log A)$, our 60 algorithms only incur an additional logarithmic overhead to account for robustness in the transition 61 probabilities. This computational complexity compares favorably with the larger time complexity of 62 a recent first-order solution scheme for KL divergence-constrained s-rectangular RMDPs (which we 63 will elaborate on later in the paper) as well as a minimum complexity of $\mathcal{O}(S^{4.5} \cdot A)$ for the naïve 64 solution with state-of-the-art interior-point algorithms. Our framework is general enough to readily 65 accommodate for ϕ -divergences that have not been studied previously in the context of s-rectangular 66 ambiguity sets, such as the Burg entropy and the χ^2 -distance. For other ϕ -divergences, such as the 67 68 L_1 -norm, our framework results in the same complexity at substantially simplified proofs.

The algorithms developed in this paper speed up the computation of robust Bellman updates and so 69 they can be used in combination with a variety of RMDP solution schemes. In particular, they can be 70 used to accelerate the standard robust value iteration, policy iteration, modified policy iteration [23] 71 and partial policy iteration [20]. They can also be combined with a first order gradient method [15] 72 that has been introduced recently. In addition, fast algorithms for computing the Bellman operator also 73 play a crucial role when scaling robust algorithms to value function approximation [52], model-free 74 reinforcement learning [34, 44], and robust policy gradients [51]. In this paper, we focus on the 75 model-based setting, which is currently under active study [25, 30, 32] and has many important real-76 life applications [13, 21, 60]; moreover, it also serves as an important building block to constructing 77 model-free algorithms. While this paper focuses on the s-rectangular ambiguity sets, the proposed 78 algorithms in this paper can also be applied to the case of (s, a)-rectangular ambiguity sets. 79

⁸⁰ The remainder of the paper proceeds as follows. Section 2 reviews relevant prior work and Section 3 describes our basic RMDP setting. Then, Section 4 shows how the robust Bellman operator for a large class of ambiguity sets can be reduced to a sequence of structured projections onto a simplex. We describe novel algorithms for efficiently computing the simplex projections for several ϕ -divergences in Section 5. Finally, Section 6 presents experimental results that compare the runtime of our algorithms with general conic solvers as well as a recent first-order optimization algorithm [15].

86 Notation. We denote by e the vector of all ones, whose context determines its dimension. We 87 refer to the probability simplex in \mathbb{R}^n by $\Delta_n = \{ p \in \mathbb{R}^n_+ : e^\top p = 1 \}$. For $x \in \mathbb{R}^n$, we let 88 min $\{x\} = \min\{x_i : i = 1, ..., n\}$ (similar for the maximum operator), and we define $[x]_+ \in \mathbb{R}^n_+$ 90 component-wise as $([x]_+)_i = \max\{x_i, 0\}, i = 1, ..., n$. We refer to the conjugate of a function 90 $f : \mathbb{R}^n \to \mathbb{R}$ by $f^*(y) = \sup\{y^\top x - f(x) : x \in \mathbb{R}^n\}$. Random variables are indicated by a tilde.

91 2 Related Work

While RMDPs have been studied since the seventies [47], they have witnessed significant recent
interest due to their widespread adoption in applications ranging from assortment optimization [45],
medical decision-making [13, 64] and hospital operations management [17], production planning [60]
and energy systems [21] to model predictive control [11], aircraft collision avoidance [24], wireless
communications [59] and the robustification against approximation errors in aggregated MDPs [38].

Efficient implementations of the robust value iteration have been first proposed by [12, 22, 33] 97 for RMDPs with (s, a)-rectangular ambiguity sets, where the worst transition probabilities are 98 considered separately for each state and action. The authors study ambiguity sets that bound the 99 distance of the transition probabilities to some nominal distribution in terms of finite scenarios, 100 interval matrix bounds, ellipsoids, the relative entropy, the KL divergence and maximum a posteriori 101 models. Subsequently, similar methods have been developed by [59] for interval matrix bounds as 102 well as likelihood uncertainty models, by [38] for 1-norm ambiguity sets as well as by [64] for interval 103 matrix bounds intersected with a budget constraint. All of these contributions have in common that 104 they focus on (s, a)-rectangular ambiguity sets where the existence of optimal deterministic policies 105 is guaranteed, and it is not clear how they could be extended to the more general class of s-rectangular 106 ambiguity sets where all optimal policies may be randomized. 107

In contrast to (s, a)-rectangular ambiguity sets, s-rectangular ambiguity sets restrict the conservatism 108 among transition probabilities corresponding to different actions in the same state, which tends to 109 lead to a superior performance in data-driven settings. [58] solve the subproblems arising in the 110 robust value iteration of an s-rectangular RMDP as linear or conic optimization problems using 111 commercial off-the-shelf solvers. Despite their polynomial-time complexity, general-purpose solvers 112 cannot exploit the structure present in these subproblems, which renders them suitable primarily 113 for small problem instances. More efficient tailored solution methods for s-rectangular RMDPs 114 have subsequently been developed by [5, 19, 20]. [19] develop a homotopy continuation method for 115 RMDPs with (s, a)-rectangular and s-rectangular weighted 1-norm ambiguity sets, while [5] adapt 116 the algorithm of [19] to unweighted ∞ -norm ambiguity sets. [20] embed the algorithms of [19] in a 117 partial policy iteration, which generalizes the robust modified policy iteration proposed by [23] for 118 (s, a)-rectangular RMDPs to s-rectangular RMDPs. 119

While the present paper focuses on the robust value iteration for ease of exposition, we note that our 120 algorithms can also be combined with the partial policy iteration of [20] to obtain further speedups. 121 [9] establish a relationship between s-rectangular RMDPs and twice regularized MDPs, which they 122 subsequently use to propose efficient Bellman updates for a modified policy iteration. While their 123 approach can solve RMDPs in almost the same time as a classical non-robust MDPs, the obtained 124 policies can be conservative as the worst-case transition probabilities are not restricted to reside in a 125 probability simplex and, therefore, may be negative and/or add up to more or less than 1. Finally, 126 [15] propose a first-order framework for RMDPs with s-rectangular KL and spherical ambiguity sets 127 128 that interleaves primal-dual first-order updates with approximate value iteration steps. The authors show that their algorithms outperform a robust value iteration that solves the emerging subproblems 129 using state-of-the-art commercial solvers. We compare our solution method for KL ambiguity sets 130 with the approach proposed by [15] in terms of its theoretical complexity and numerical runtimes. 131

While this paper exclusively studies s-rectangular uncertainty sets, alternative generalizations of (s, a)-132 rectangular ambiguity sets have been proposed in the literature as well. For example, [29] consider 133 k-rectangular ambiguity sets where the transition probabilities of different states can be coupled, [14] 134 study factor ambiguity model ambiguity sets where the transition probabilities depend on a small 135 number of underlying factors, and [53] construct ambiguity sets that bound marginal moments of 136 state-action features defined over entire MDP trajectories. Other than model-based settings, there 137 is also an interesting line of research on robust reinforcement learning, such as least squares policy 138 iteration [34], analysis on sample complexity [35], robust Q-learning algorithm and robust TDC 139 algorithm [44, 55], and robust policy gradient [56]. We also note the papers [7, 16, 62] which study 140 the related problem of *distributionally* robust MDPs whose transition probabilities are themselves 141 regarded as random objects that are drawn from distributions which are only partially known. The 142 connections between RMDPs and multi-stage stochastic programs as well as distributionally robust 143 problems are explored further by [46, 48, 49]. 144

145 **3** Preliminaries

Robust MDPs We study RMDPs with a finite state space $S = \{1, ..., S\}$ and a finite action space 146 $\mathcal{A} = \{1, \ldots, A\}$. We assume an infinite planning horizon, but all of our results immediately extend 147 to a finite time horizon. Without loss of generality, we assume that every action $a \in \mathcal{A}$ is admissible 148 in every state $s \in S$. The RMDP starts in a random initial state \tilde{s}_0 that follows the known probability 149 distribution p^0 from the probability simplex Δ_S in \mathbb{R}^S . If action $a \in \mathcal{A}$ is taken in state $s \in S$, then 150 the RMDP transitions randomly to the next state according to the conditional probability distribution 151 $p_{sa} \in \Delta_S$. We condense the transition probabilities p_{sa} to the tensor $p \in (\Delta_S)^{S \times A}$. The transition probabilities are only known to reside in a non-empty, compact ambiguity set $\mathcal{P} \subseteq (\Delta_S)^{S \times A}$. For 152 153 a transition from state $s \in S$ to state $s' \in S$ under action $a \in A$, the decision maker receives an 154 expected reward of $r_{sas'} \in \mathbb{R}_+$. As with the transition probabilities, we condense these rewards to the tensor $r \in \mathbb{R}^{S \times A \times S}_+$. Without loss of generality, we assume that all rewards are non-negative. 155 156

We denote by $\Pi = (\Delta_A)^S$ the set of all stationary (*i.e.*, time-independent) randomized policies. A policy $\pi \in \Pi$ takes action $a \in \mathcal{A}$ in state $s \in \mathcal{S}$ with probability π_{sa} . The transition probabilities $p \in \mathcal{P}$ and the policy $\pi \in \Pi$ induce a stochastic process $\{(\tilde{s}_t, \tilde{a}_t)\}_{t=0}^{\infty}$ on the space $(\mathcal{S} \times \mathcal{A})^{\infty}$ of sample paths. We refer by $\mathbb{E}^{p,\pi}$ to expectations with respect to this process. The decision maker is risk-neutral but ambiguity-averse and wishes to maximize the worst-case expected total reward under a discount factor $\lambda \in (0, 1)$,

$$\max_{\boldsymbol{\pi}\in\Pi} \min_{\boldsymbol{p}\in\mathcal{P}} \mathbb{E}^{\boldsymbol{p},\boldsymbol{\pi}} \left[\sum_{t=0}^{\infty} \lambda^t \cdot r_{\tilde{s}_t, \tilde{a}_t, \tilde{s}_{t+1}} \mid \tilde{s}_0 \sim \boldsymbol{p}^0 \right].$$
(1)

Note that the maximum and minimum in (1) are both attained by the Weierstrass theorem since Π and \mathcal{P} are non-empty and compact, while the objective function is finite since $\lambda < 1$.

Rectangular Ambiguity Sets For general ambiguity sets \mathcal{P} , evaluating the inner minimization in (1) is NP-hard even if the policy $\pi \in \Pi$ is fixed [58]. For these reasons, much of the research on RMDPs and their applications has focused on rectangular ambiguity sets. Among the most general rectangular ambiguity sets are the *s*-rectangular ambiguity sets \mathcal{P} satisfying

$$\mathcal{P} = \left\{ \boldsymbol{p} \in (\Delta_S)^{S \times A} : \, \boldsymbol{p}_s \in \mathcal{P}_s \ \forall s \in \mathcal{S} \right\}, \quad \text{where} \quad \mathcal{P}_s \subseteq (\Delta_S)^A, \, s \in \mathcal{S},$$

see [26, 58, 61, 63]. In contrast to the simpler class of (s, a)-rectangular ambiguity sets, *s*-rectangular ambiguity sets restrict the choice of transition probabilities p_{s1}, \ldots, p_{sA} corresponding to different actions *a* applied in the same state *s*. This limits the conservatism of the resulting RMDP (1) and typically leads to a better performance of the optimal policy [58]. Although Bellman's optimality principle extends to *s*-rectangular RMDPs and there is always an optimal stationary policy, all optimal policies of an *s*-rectangular RMDP may be randomized.

¹⁷⁵ We study a new general class of *s*-rectangular ambiguity sets that can be expressed as

$$\mathcal{P}_{s} = \left\{ \boldsymbol{p}_{s} \in (\Delta_{S})^{A} : \sum_{a \in \mathcal{A}} d_{a}(\boldsymbol{p}_{sa}, \overline{\boldsymbol{p}}_{sa}) \leq \kappa \right\},$$
(2)

where $\kappa \in \mathbb{R}_+$ is the *uncertainty budget* and the distance functions $d_a(\mathbf{p}_{sa}, \overline{\mathbf{p}}_{sa}), a \in \mathcal{A}$, are ϕ -divergences (also known as *f*-divergences) satisfying

$$d_a(\boldsymbol{p}_{sa}, \overline{\boldsymbol{p}}_{sa}) = \sum_{s' \in \mathcal{S}} \overline{p}_{sas'} \phi\left(\frac{p_{sas'}}{\overline{p}_{sas'}}\right).$$

Here, $\phi: \mathbb{R}_+ \to \mathbb{R}_+$ is a convex function satisfying $\phi(1) = 0$. Intuitively, a ϕ -divergence measures the distance between two probability distributions. With an appropriate choice of ϕ , it generalizes other metrics including the KL divergence, the Burg entropy, L_1 - and L_2 -norms and others [4, 6]. Table 1 reports some popular ϕ -divergences that we study in this paper. Note that the variation distance coincides with the L_1 -based *s*-rectangular ambiguity sets studied in earlier work [19, 20]. Note that although ϕ is set to be the same for different state-action pairs, the proposed approach also work for the general case of $d_a(\mathbf{p}_{sa}, \overline{\mathbf{p}}_{sa}) = \sum_{s' \in S} \overline{p}_{sas'} \phi_{sas'}(p_{sas'})$.

Divergence	$ d_a(oldsymbol{p}_{sa}, \overline{oldsymbol{p}}_{sa}) $	$\phi(t)$	Complexity of \mathfrak{J}	State-of-the-Art
Kullback-Leibler	$\sum_{s'} p_{sas'} \log\left(\frac{p_{sas'}}{\overline{p}_{sas'}}\right)$	$t\log t - t + 1$	$\mathcal{O}(S^2 \cdot A \log A)$	$\mathcal{O}(\ell^2 \cdot S^2 \cdot A)$
Burg Entropy	$\sum_{s'} \overline{p}_{sas'} \log\left(\frac{\overline{p}_{sas'}}{p_{sas'}}\right)$	$-\log t + t - 1$	$\mathcal{O}(S^2 \cdot A \log A)$	no poly-time guarantee
Variation Distance	$\sum_{s'} p_{sas'} - \overline{p}_{sas'} '$	t-1	$\mathcal{O}(S^2 \log S \cdot A)$	$\mathcal{O}(S^2 \log S \cdot A)$
χ^2 -Distance	$\sum_{s'} \frac{(p_{sas'} - \overline{p}_{sas'})^2}{\overline{p}_{sas'}}$	$(t - 1)^2$	$\mathcal{O}(S^2 \log S \cdot A)$	$\mathcal{O}(S^{4.5} \cdot A)$

Table 1: Summary of the ϕ -divergences studied in this paper, together with the complexity of our robust Bellman operator \mathfrak{J} (applied across all states $s \in S$) as well as the best known results from the literature. The complexity estimates omit constants and tolerances that are reported in Section 5 of the paper. ' ℓ ', where present, refers to the number of Bellman iterations conducted so far.

Robust Value Iteration A standard approach for computing the optimal value and the optimal policy of an RMDP (1) is the robust value iteration [22, 33, 26, 58]: Starting with an initial estimate $v^0 \in \mathbb{R}^S$ of the state-wise optimal value to-go, we conduct robust Bellman iterations of the form $v^{t+1} \leftarrow \mathfrak{J}(v^t), t = 0, 1, ...$, where the robust Bellman operator \mathfrak{J} is defined component-wise as

$$[\mathfrak{J}(\boldsymbol{v})]_{s} = \max_{\boldsymbol{\pi}_{s} \in \Delta_{A}} \min_{\boldsymbol{p}_{s} \in \mathcal{P}_{s}} \sum_{a \in \mathcal{A}} \boldsymbol{\pi}_{sa} \cdot \boldsymbol{p}_{sa}^{\top} (\boldsymbol{r}_{sa} + \lambda \boldsymbol{v}) \quad \forall s \in \mathcal{S}.$$
(3)

This yields the optimal value $p^{0^{\top}}v^{\star}$, where the limit $v^{\star} = \lim_{t\to\infty} v^t$ is approached component-wise at a geometric rate. The optimal policy $\pi^{\star} \in \Pi$, finally, is recovered state-wise via

$$\pi_s^{\star} \in \operatorname*{arg\,max}_{\pi_s \in \Delta_A} \min_{p_s \in \mathcal{P}_s} \sum_{a \in \mathcal{A}} \pi_{sa} \cdot p_{sa}^{\top} (r_{sa} + \lambda v^{\star}) \quad \forall s \in \mathcal{S}.$$

191 4 Robust Bellman Updates via Simplex Projections

¹⁹² In this section, we show that the robust Bellman operator \mathfrak{J} reduces to a generalized projection ¹⁹³ problem. This reduction is important because it underlies our fast algorithms for computing \mathfrak{J} .

194 At the core of the robust value iteration is the solution of the max-min problem (3). By applying

min-max theorem, the optimal value in right-hand side of (3) could be reduced to be a structural

problem which could be solved via bisection on its objective value. For any given β in the bisection

¹⁹⁷ method, we check whether feasible p_s exists via solving the following generalized d_a -projection of ¹⁹⁸ the nominal transition probabilities \overline{p}_{sa} :

$$\mathfrak{P}(\overline{\boldsymbol{p}}_{sa}; \boldsymbol{b}, \beta) = \begin{bmatrix} \text{minimize} & d_a(\boldsymbol{p}_{sa}, \overline{\boldsymbol{p}}_{sa}) \\ \text{subject to} & \boldsymbol{b}^\top \boldsymbol{p}_{sa} \le \beta \\ & \boldsymbol{p}_{sa} \in \Delta_S \end{bmatrix}.$$
(4)

Here, $p_{sa} \in \Delta_S$ are the decision variables and $\overline{p}_{sa} \in \Delta_S$, $b \in \mathbb{R}^S_+$ and $\beta \in \mathbb{R}_+$ are parameters. Thus, the robust Bellman operator \mathfrak{J} could be computed efficiently if the above generalized d_a -projection (4) could be computed efficiently, and more details information on the bisection method could be found in the Appendix. Note that problem (4) is infeasible if and only if $\min\{b\} > \beta$. Moreover, problem (4) is trivially solved by \overline{p}_{sa} with an optimal objective value of 0 whenever $b^{\top} \overline{p}_{sa} \leq \beta$. To avoid these trivial cases, we assume throughout the paper that $\min\{b\} \leq \beta$ and $b^{\top} \overline{p}_{sa} > \beta$. We illustrate the feasible region and optimal solution to problem (4) for different ϕ -divergences in Figure 1.

Our generalized d_a -projection (4) relates to the rich literature on projections onto simplices, which we review in the next section. In fact, our algorithms in the next section solve a variant of the simplex projection problem that is restricted by an additional inequality constraint. We therefore believe that our algorithms may find additional applications outside the RMDP literature.

In the following, we say that for a given estimate $v^t \in \mathbb{R}^S$ of the optimal value function, the robust Bellman iteration (3) is solved to ϵ -accuracy by any $v^{t+1} \in \mathbb{R}^S$ satisfying $\|v^{t+1} - \mathfrak{J}(v^t)\|_{\infty} \leq \epsilon$. We seek ϵ -optimal solutions because our ambiguity sets are nonlinear and hence the exact Bellman



Figure 1: The generalized d_a -projection problem (4) in S = 3 dimensions (a) and two-dimensional projections for the variation distance (b), the χ^2 -distance (c) and the KL divergence (d). The gray shaded areas represent the probability simplex Δ_S , the red dashed lines show the boundary of the intersection of the halfspace $b^{\top} p_{sa} \leq \beta$ with the probability simplex, and the white shapes illustrate contour lines centered at the nominal transition probabilities \overline{p}_{sa} .

iterate $\mathfrak{J}(\boldsymbol{v}^t)$ may be irrational even if \boldsymbol{v}^t is rational. To simplify the exposition, we define $\overline{R} = 14$ $[1-\lambda]^{-1} \cdot \max\{r_{sas'}: s, s' \in \mathcal{S}, a \in \mathcal{A}\}$ as an upper bound on all $[\mathfrak{J}(\boldsymbol{v})]_s, \boldsymbol{v} \leq \boldsymbol{v}^*$ and $s \in \mathcal{S}$.

For divergence-based ambiguity sets, the projection problem (4) is generically nonlinear and can hence not be expected to be solved to exact optimality. To account for this additional complication, we say that for a given $\overline{p}_{sa} \in \Delta_S$, $\mathbf{b} \in \mathbb{R}^S_+$ and $\beta \in \mathbb{R}_+$, the generalized d_a -projection $\mathfrak{P}(\overline{p}_{sa}; \mathbf{b}, \beta)$ is solved to δ -accuracy by any pair $(\underline{d}, \overline{d}) \in \mathbb{R}^2$ satisfying $\mathfrak{P}(\overline{p}_{sa}; \mathbf{b}, \beta) \in [\underline{d}, \overline{d}]$ and $\overline{d} - \underline{d} \leq \delta$.

Theorem 1. Assume that the generalized d_a -projection (4) can be computed to any accuracy $\delta > 0$ in time $\mathcal{O}(h(\delta))$. Then the robust Bellman iteration (3) can be computed to any accuracy $\epsilon > 0$ in time $\mathcal{O}(AS \cdot h(\epsilon \kappa / [2A\overline{R} + A\epsilon]) \cdot \log[\overline{R}/\epsilon])$.

Theorem 1 reduces the evaluation of the robust Bellman iterator \mathfrak{J} , which involves the solution of a 222 max-min optimization problem over an s-rectangular ambiguity set that couples all actions $a \in \mathcal{A}$, to 223 a sequence of much simpler and highly structured projection problems that are no longer coupled 224 225 across different actions $a \in A$. The next section describes efficient solution schemes for the projection problem (4) in the context of several ϕ -divergence ambiguity sets. The runtimes of these solution 226 schemes are summarized in Table 1. Note that the evaluation of a non-robust Bellman operator 227 requires a runtime of $\mathcal{O}(S^2 \cdot \log A)$, which implies that our algorithms only incur an additional 228 logarithmic overhead to account for robustness in the transition probabilities. 229

230 5 Fast Projections on ϕ -Divergence Simplices

We next describe fast algorithms for computing generalized projections onto the probability simplex. 231 Combined with the results from Section 4, these algorithms can be used to efficiently compute 232 the robust Bellman operator. Note that some ϕ -divergences, such as the KL divergence and the 233 χ^2 -distance, imply that if $\overline{p}_{sas'} = 0$ for some $s, s' \in S$ and $a \in A$, then $p_{sas'} = 0$ for all $p_{sa} \in \Delta_S$ 234 with $d_a(p_{sa}, \overline{p}_{sa}) < \infty$, and thus we can remove indices s' with $\overline{p}_{sas'} = 0$. For other ϕ -divergences, 235 such as the Burg entropy and the variation distance, one can readily verify that our results remain 236 valid no matter whether $\overline{p}_{sa} > 0$ or not, but the formulations and proofs require additional case 237 distinctions and/or limit arguments. To simplify the exposition, we therefore assume that $\overline{p}_{sa} > 0$. 238

Proposition 1. For the distance function $d_a(\mathbf{p}_{sa}, \overline{\mathbf{p}}_{sa}) = \sum_{s' \in S} \overline{p}_{sas'} \cdot \phi\left(\frac{p_{sas'}}{\overline{p}_{sas'}}\right)$, the optimal value of the projection problem (4) equals the optimal value of the bivariate convex problem

maximize
$$-\beta \alpha + \zeta - \sum_{s' \in S} \overline{p}_{sas'} \phi^{\star}(-\alpha b_{s'} + \zeta)$$

subject to $\alpha \in \mathbb{R}_+, \ \zeta \in \mathbb{R}.$ (5)

Proposition 1 reduces the S-dimensional projection problem (4) to a two-dimensional optimization

problem over the dual variables α and ζ . In the following, we show that for the ϕ -divergences from Table 1, problem (5) can be further simplified to univariate convex optimization problems that can be

solved efficiently via bisection, binary search or sorting.

245 5.1 Kullback-Leibler Divergence

We first show that for the KL divergence $\phi(t) = t \log t - t + 1$, the reduced projection problem (5) can be further simplified to a univariate convex optimization problem. The radii of this type of ambiguity sets can be selected via statistical bounds [27].

Proposition 2. For the KL divergence $\phi(t) = t \log t - t + 1$, the optimal value of the projection problem (4) equals the optimal value of the univariate convex problem

$$\underset{\alpha \in \mathbb{R}_{+}}{\operatorname{maximize}} \quad -\beta\alpha - \log\left(\sum_{s' \in \mathcal{S}} \overline{p}_{sas'} \cdot e^{-\alpha b_{s'}}\right).$$
(6)

We next show that the univariate optimization problem (2) admits an efficient solution via bisection. **Theorem 2.** If $\beta \ge \min\{b\} + \omega$ for some $\omega > 0$, then the projection problem (4) can be solved to any δ -accuracy in time $\mathcal{O}(S \cdot \log[\max\{b\} \cdot \log(\min\{\overline{p}\}^{-1})/(\delta\omega)])$.

Note that the projection problem (4) is infeasible whenever $\beta < \min\{b\}$. The condition in the 254 statement of Theorem 2 can thus be interpreted as a strict feasibility requirement. It is worth 255 contrasting the result of Theorem 2 with the solution of the projection problem (4) as an exponential 256 cone program. The latter would result in a *practical* complexity of $\mathcal{O}(S^3)$, assuming that—which is 257 often observed in practice-the number of iterations of the employed interior-point solver does not 258 grow with the problem dimensions. A theoretically guaranteed complexity, on the other hand, does 259 not seem to be available at present as the commercial state-of-the-art solvers for exponential conic 260 programs are not proven to terminate in polynomial time. 261

Corollary 1. The robust Bellman iteration (3) over a KL divergence ambiguity set can be computed to any accuracy $\epsilon > 0$ in time $\mathcal{O}(S^2 \cdot A \log A \cdot \log[\overline{R}^2 \cdot \log(\min\{\overline{p}\}^{-1})/(\epsilon^2 \kappa)] \cdot \log[\overline{R}/\epsilon])$.

[15] propose a first-order framework for RMDPs over s-rectangular KL divergence ambiguity sets 264 whose robust Bellman update enjoys a complexity of $\mathcal{O}(\ell^2 \cdot S^2 \cdot A \cdot \log(\epsilon^{-1}))$, where ℓ is the iteration 265 number. A careful analysis results in an overall convergence rate for the optimal MDP policy of 266 $\mathcal{O}(S^3 \cdot A^2 \cdot \epsilon^{-1} \log[\epsilon^{-1}])$. In contrast, the convergence rate of our robust value iteration amounts to 267 $\mathcal{O}(S^2 \cdot A \log A \cdot \log[\overline{R}^2 \cdot \log(\min\{\overline{p}\}^{-1})/(\epsilon^2 \kappa)] \cdot \log[\overline{R}/\epsilon] \cdot \log[\epsilon^{-1}]).$ Treating the problem parameters $\overline{R}, \overline{p}$ and κ as constants, our convergence rate simplifies to $\mathcal{O}(S^2 \cdot A \log A \cdot \log[\epsilon^{-2}] \cdot \log^2[\epsilon^{-1}]),$ 268 269 which compares favourably against the convergence rate of the first-order scheme. Our numerical 270 results in Section 6 show that this theoretical difference appears to carry over to a favourable empirical 271 performance on test instances as well. 272

We finally note the related work [1], which optimizes a linear function over the intersection of a probability simplex with a constraint on the KL divergence to a nominal distribution. While one could in principle modify that algorithm to solve our projection problem (4), the resulting algorithm would require an additional bisection and would thus be significantly slower than ours.

277 5.2 Burg Entropy

Similar to the KL divergence, the reduced projection problem (5) can be further simplified to a univariate convex optimization problem for the Burg entropy $\phi(t) = -\log t + t - 1$. The radii of this type of ambiguity sets can be selected via statistical bounds [27].

Proposition 3. For the Burg entropy $\phi(t) = -\log t + t - 1$, if $\beta > \min\{b\}$, then the optimal value of the projection problem (4) equals the optimal value of the univariate convex problem

$$\underset{\alpha \in [0,1]}{\operatorname{maximize}} \quad \sum_{s' \in \mathcal{S}} \overline{p}_{sas'} \cdot \log\left(1 + \alpha \frac{b_{s'} - \beta}{\beta - \min\{\boldsymbol{b}\}}\right).$$
(7)

283 Similar to the KL divergence, the univariate optimization problem (7) can be solved efficiently.

Theorem 3. If $\beta \ge \min\{\mathbf{b}\} + \omega$ for some $\omega > 0$, then the projection problem (4) can be solved to any δ -accuracy in time $\mathcal{O}(S \cdot \log[\max\{\mathbf{b}\}/(\delta\omega)])$.

As with the KL divergence, the projection problem (4) corresponding to the Burg entropy can be achieved in a prostical complexity of $\mathcal{O}(\mathcal{C}^3)$ as an exponential one program whereas we are not away

- of any state-of-the-art solvers equipped with theoretical guarantees. To our best knowledge, RMDPs with *s*-rectangular Burg entropy ambiguity sets have not been studied previously in the literature.
- **Corollary 2.** The robust Bellman iteration (3) over a Burg entropy ambiguity set can be computed to any accuracy $\epsilon > 0$ in time $\mathcal{O}(S^2 \cdot A \log A \cdot \log[\overline{R}^2/(\epsilon^2 \kappa)] \cdot \log[\overline{R}/\epsilon])$.
- Similar to the previous subsection, we note that the related paper [1] optimizes a linear function over the intersection of a probability simplex with a bound on the Burg entropy to a nominal distribution.
- While that algorithm could in principle be employed to solve our projection problem (4), the resulting solution scheme would not be competitive due to the inclusion of an additional bisection.
- 295 solution scheme would not be competitive due to the inclusion of an additional disection.

296 5.3 Variation Distance

- We first provide an equivalent univariate optimization problem for the reduced projection problem (5) corresponding to the variation distance $\phi(t) = |t - 1|$. The radii of this type of ambiguity sets can be selected via statistical bounds [57].
- **Proposition 4.** For the variation distance $\phi(t) = |t 1|$, the optimal value of the projection problem (4) equals the optimal value of the univariate convex problem

$$\underset{\alpha \in \mathbb{R}_{+}}{\text{maximize}} \ 2 + \alpha(\min\{\boldsymbol{b}\} - \beta) - \sum_{s' \in \mathcal{S}} \overline{p}_{sas'} \cdot [2 + \alpha \cdot (\min\{\boldsymbol{b}\} - b_{s'})]_{+}.$$
(8)

302 Once more, the univariate optimization problem (8) admits an efficient solution.

Theorem 4. The projection problem (4) can be solved exactly in time $\mathcal{O}(S \log S)$.

Note that in contrast to the previous results, Theorem 4 employs a binary search and thus offers an *exact* solution to the projection problem (4). Our result of Theorem 4 matches the complexity of the homotopy continuation method proposed by [20]. The correctness and runtime of their algorithm, however, relies on lengthy ad hoc arguments, whereas Theorem 4 relies on the groundwork laid by Theorem 1 and Proposition 1. Problem (4) can also be solved as a linear program with a practical complexity of $\mathcal{O}(S^3)$ and a theoretical complexity of $\mathcal{O}(S^{3.5})$.

Corollary 3. The robust Bellman iteration (3) over a variation distance ambiguity set can be computed to any accuracy $\epsilon > 0$ in time $\mathcal{O}(S^2 \log S \cdot A \cdot \log[\overline{R}/\epsilon])$.

study the related problem of optimizing a linear function over the intersection of a probability simplex with an unweighted 1-norm constraint, and they identify structural properties of the optimal solutions. Since the linear function and the norm constraint are in different places of the optimization problem, however, their findings are not directly applicable to our setting.

316 **5.4** χ^2 -Distance

In contrast to the previous subsections, we directly solve the bivariate problem (5) for the χ^2 -distance $\phi(t) = (t-1)^2$ without first formulating an associated univariate optimization problem. The radii this type of ambiguity sets can be selected via statistical bounds [33, 57].

Theorem 5. For the χ^2 -distance $\phi(t) = (t-1)^2$, the optimal value of the projection problem (4) can be computed exactly in time $\mathcal{O}(S \log S)$.

Theorem 5 splits the bivariate piecewise quadratic optimization problem (5) corresponding to the χ^2 -distance into S + 1 bivariate quadratic problems by sorting the components of **b**. Each of these S + 1 problems can be reduced to the solution of 3 univariate quadratic problems that themselves admit analytical solutions.

Corollary 4. The robust Bellman iteration (3) over a χ^2 -distance ambiguity set can be computed to any accuracy $\epsilon > 0$ in time $\mathcal{O}(S^2 \log S \cdot A \cdot \log[\overline{R}/\epsilon])$.

The projection problem (4) for the χ^2 -distance ambiguity set can be solved as a quadratic program with a practical complexity of $\mathcal{O}(S^3)$ as well as a theoretical complexity of $\mathcal{O}(S^{3.5})$.

- ³³⁰ The first-order framework of [15] also applies to RMDPs over *s*-rectangular spherical uncertainty
- sets. In that case, the robust Bellman update enjoys a complexity of $\mathcal{O}(\ell^2 \cdot S^2 \cdot A \cdot \log^2(\epsilon^{-1}))$, where

S	MOSEK	fast	MOSEK/fast		S = A	MOSEK	fast	MOSEK/fast
20	1.00	0.01	175.35	-	20	12.98	1.06	12.21
100	7.53	0.02	317.80		100	637.78	25.25	25.28
400	17.87	0.09	190.95		400	24,308.16	343.37	70.79
1,000	49.23	0.24	208.20		600	47,473.61	731.17	64.93
4,000	235.43	0.94	249.18		700	63,318.00	1,084.65	58.38

Table 2: Comparison of our algorithms ('fast') vs. MOSEK for the projection problem (left) and the Bellman update (right) on KL-divergence constrained ambiguity sets. Runtimes are reported in ms.

S = A	f-o (3 its)	f-o (5 its)	fast	f-o/fast (3 its)	f-o/fast (5 its)
20	9.12	25.25	1.06	8.58	23.75
100	183.34	508.83	25.25	7.26	20.15
400	2,821.52	7,833.65	343.37	8.21	22.81
600	6,434.55	17,828.39	731.17	8.80	24.38
700	8,523.80	23,702.00	1,084.65	7.86	21.85

Table 3: Comparison of our algorithms ('fast') vs. the first-order method of [15] (after $\ell = 3, 5$ its.) for the Bellman update on KL-divergence constrained ambiguity sets. Runtimes are reported in ms.

³³² ℓ is the iteration number. A careful analysis results in an overall convergence rate for the optimal ³³³ MDP policy of $\mathcal{O}(S^3 \log S \cdot A^2 \cdot \epsilon^{-1} \log[\epsilon^{-1}])$. In contrast, the convergence rate of our robust value ³³⁴ iteration amounts to $\mathcal{O}(S^2 \log S \cdot A \cdot \log[\overline{R}/\epsilon] \cdot \log[\epsilon^{-1}])$. Treating the parameter \overline{R} as a constant, our ³³⁵ convergence rate simplifies to $\mathcal{O}(S^2 \log S \cdot A \cdot \log[\overline{R}/\epsilon])$, which compares favourably against the ³³⁶ convergence rate of [15]. We remark, however, that the spherical ambiguity sets of [15] differ from ³³⁷ the χ^2 -distance ambiguity sets studied here, and as such the two methods are not directly comparable. ³³⁸ We also note that our χ^2 -distance ambiguity sets enjoy a strong statistical justification [4, 6].

Computing unweighted 2-norm projections of points onto S-dimensional probability simplices has 339 manifold applications in image processing, finance, optimization and machine learning [1, 8]. [31] 340 proposes one of the earliest algorithms that computes this projection in time $\mathcal{O}(S^2)$ by iteratively 341 reducing the dimension of the problem using Lagrange multipliers. The minimum complexity of 342 $\mathcal{O}(S)$ is achieved, among others, by [28] through a linear-time median-finding algorithm and by [37] 343 through a filtered bucket-clustering method. Note, however, that these algorithms do not account 344 for the weights and the additional inequality constraint present in our generalized projection (4). 345 The unweighted 2-norm projection of a point onto the intersection of the S-dimensional probability 346 simplex with an axis-parallel hypercube is computed by [54] through a sorting-based method and 347 by [2] through Newton's method, respectively. [39] optimize a linear function over the intersection 348 of a probability simplex with an unweighted 2-norm constraint through an iterative dimension 349 reduction scheme. [1], finally, study algorithms that optimize linear functions over the intersection of 350 a probability simplex and a bound on the unweighted 2-norm distance to a nominal distribution. 351

352 6 Numerical Results

We compare our fast suite of algorithms with the state-of-the-art solver MOSEK 9.3 [3] (commercial) and the first-order method of [15]. Tables 2–3 report average computation times over 50 randomly generated test instances for the KL-divergence case, and show that the proposed algorithms outperforms other methods. Similar experimental results for the χ^2 -distance is provided in the Appendix, which provides the details of all the experiments.

358 7 Conclusion

We consider the robust MDPs with *s*-rectangular ϕ -divergence ambiguity sets. We develop efficient algorithms for computing the robust Bellman updates for several important special cases of this ambiguity set. Our experimental results indicate that the proposed algorithms outperform MOSEK.

³⁶² Future work should address extensions to the developments of scalable model-free algorithms.

363 References

- [1] L. Adam and V. Mácha. Projections onto the canonical simplex with additional linear inequalities.
 Optimization Methods & Software, Available online first, 2020.
- [2] M. S. Ang, J. Ma, N. Liu, K. Huang, and Y. Wang. Fast projection onto the capped simplex
 with applications to sparse regression in bioinformatics. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [3] MOSEK ApS. MOSEK Fusion API for C++ 9.3.20, 2019. URL https://docs.mosek.com/
 latest/cxxfusion/index.html.
- [4] G. Bayraksan and D. K. Love. Data-driven stochastic programming using phi-divergences. In
 D. M. Aleman and A. C. Thiele, editors, *INFORMS TutORials in Operations Research*, pages
 1–19. 2015.
- [5] B. Behzadian, M. Petrik, and C. P. Ho. Fast algorithms for l_{∞} -constrained s-rectangular robust MDPs. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [6] A. Ben-Tal, D. den Hertog, A. de Waegenaere, B. Melenberg, and G. Rennen. Robust solutions
 of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):
 341–357, 2013.
- [7] Z. Chen, P. Yu, and W. B. Haskell. Distributionally robust optimization for sequential decision making. *Optimization*, 68(12):2397–2426, 2019.
- [8] L. Condat. Fast projection onto the simplex and the l_1 ball. *Mathematical Programming*, 158 (1–2):575–585, 2016.
- [9] E. Derman, M. Geist, and S. Mannor. Twice regularized MDPs and the equivalence between ro bustness and regularization. In *Advances in Neural Information Processing Systems*, volume 35,
 pages (Pre–Proceedings), 2021.
- [10] S. Diamond and S. Boyd. CVXPY: A Python-embedded modeling language for convex
 optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [11] M. Diehl and J. Bjornberg. Robust dynamic programming for min-max model predictive control of constrained uncertain systems. *IEEE Transactions on Automatic Control*, 49(12):2253–2257, 2004.
- [12] R. Givan, S. Leach, and T. Dean. Bounded-parameter Markov decision processes. *Artificial Intelligence*, 122(1):71–109, 2000.
- J. Goh, M. Bayati, S. A. Zenios, S. Singh, and D. Moore. Data uncertainty in Markov chains:
 Application to cost-effectiveness analyses of medical innovations. *Operations Research*, 66(3):
 697–715, 2018.
- [14] V. Goyal and J. Grand-Clément. Robust Markov decision process: Beyond rectangularity.
 Available on arXiv, 2018.
- [15] J. Grand-Clément and C. Kroer. Scalable first-order methods for robust MDPs. In *Proceedings* of the AAAI Conference on Artificial Intelligence, volume 35, pages 12086–12094, 2021.
- [16] J. Grand-Clément and C. Kroer. First-order methods for Wasserstein distributionally robust
 MDPs. In *Proceedings of Machine Learning Research*, volume 139, pages 2010–2019, 2021.
- [17] J. Grand-Clément, C. W. Chan, V. Goyal, and G. Escobar. Robust policies for proactive ICU
 transfers. Working Paper, 2019.
- 404 [18] Vishal Gupta. Near-optimal Bayesian ambiguity sets for distributionally robust optimization.
 405 Management Science, 65(9):4242–4260, 2019.
- [19] C. P. Ho, M. Petrik, and W. Wiesemann. Fast Bellman updates for robust MDPs. In *Proceedings* of the 35th International Conference on Machine Learning, pages 979–1988, 2018.

- [20] C. P. Ho, M. Petrik, and W. Wiesemann. Partial policy iteration for l₁-robust Markov decision
 processes. *Journal of Machine Learning Research*, 22:1–46, 2021.
- [21] Q. Huang, Q.-S. Jia, and X. Guan. Robust scheduling of EV charging load with uncertain wind
 power integration. *IEEE Transactions on Smart Grid*, 9(2):1043–1054, 2018.
- 412 [22] G. N. Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):
 413 257–280, 2005.
- [23] D. L. Kaufman and A. J. Schaefer. Robust modified policy iteration. *INFORMS Journal on Computing*, 25(3):396–410, 2013.
- [24] M. J. Kochenderfer and J. P. Chryssanthacopoulos. Robust airborne collision avoidance through
 dynamic programming. Project Report ATC-371 for the Federal Aviation Administration, 2011.
- [25] M.A.S. Kolarijani, G.F. Max, and P. Mohajerin Esfahani. Fast approximate dynamic programming for infinite-horizon markov decision processes. In *Advances in Neural Information Processing Systems*, volume 34, pages 23652–23663, 2021.
- [26] Y. Le Tallec. *Robust, Risk-Sensitive, and Data-driven Control of Markov Decision Processes.* PhD thesis, Massachusetts Institute of Technology, 2007.
- 423 [27] D. Love and G. Bayraksan. Phi-divergence constrained ambiguous stochastic programs for 424 data-driven optimization. *Technical report*, 2015.
- [28] N. Maculan and G. G. de Paula Jr. A linear-time median-finding algorithm for projecting a vector on the simplex of \mathbb{R}^n . Operations Research Letters, 8(4):219–222, 1989.
- [29] S. Mannor, O. Mebel, and H. Xu. Robust MDPs with *k*-rectangular uncertainty. *Mathematics of Operations Research*, 41(4):1484–1509, 2016.
- [30] A. Al Marjani, A. Garivier, and A. Proutiere. Navigating to the best policy in markov decision
 processes. In *Advances in Neural Information Processing Systems*, volume 34, pages 25852–
 25864, 2021.
- ⁴³² [31] C. Michelot. A finite algorithm for finding the projection of a point onto the canonical simplex ⁴³³ of \mathbb{R}^n . *Journal of Optimization Theory and Applications*, 50(1):195–200, 1986.
- [32] A. Nie, E. Brunskill, and C. Piech. Play to grade: Testing coding games as classifying markov
 decision process. In *Advances in Neural Information Processing Systems*, volume 34, pages
 1506–1518, 2021.
- [33] A. Nilim and L. El Ghaoui. Robust control of Markov decision processes with uncertain
 transition matrices. *Operations Research*, 53(5):780–798, 2005.
- [34] K. Panaganti and D. Kalathil. Robust reinforcement learning using least squares policy iteration
 with provable performance guarantees. In *Proceedings of the 38th International Conference on Machine Learning*, pages 511–520, 2021.
- K. Panaganti and D. Kalathil. Sample complexity of robust reinforcement learning with a
 generative model. In *International Conference on Artificial Intelligence and Statistics*, pages
 9582–9602. PMLR, 2022.
- [36] B. P.G. Van Parys, P. Mohajerin Esfahani, and D. Kuhn. From data to decisions: Distributionally
 robust optimization is optimal. *Management Science*, 67(6):3387–3402, 2021.
- ⁴⁴⁷ [37] G. Perez, M. Barlaud, L. Fillatre, and J.-C. Régin. A filtered bucket-clustering method for ⁴⁴⁸ projection onto the simplex and the l_1 ball. *Mathematical Programming*, 182(1–2):445–464, ⁴⁴⁹ 2020.
- [38] M. Petrik and D. Subramanian. RAAM: The benefits of robustness in approximating aggregated
 MDPs in reinforcement learning. In *Advances in Neural Information Processing Systems*,
 volume 27, pages 1979–1987, 2014.

- [39] A. Philpott, V. de Matos, and L. Kapelevich. Distributionally robust SDDP. *Computational Management Science*, 15(3–4):431–454, 2018.
- [40] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John
 Wiley & Sons, 1994.
- [41] H. Rahimian, G. Bayraksan, and T. Homem-de-Mello. Identifying effective scenarios in distributionally robust stochastic programs with total variation distance. *Mathematical Programming*, 173(1–2):393–420, 2019.
- 460 [42] R. T. Rockafellar. Convex Analysis. Princeton University Press, 1997.
- [43] G. Van Rossum and F. L. Drake. *Python 3 Reference Manual*. CreateSpace, Scotts Valley, CA,
 2009. ISBN 1441412697.
- [44] A. Roy, H. Xu, and S. Pokutta. Reinforcement learning under model mismatch. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- [45] P. Rusmevichientong and H. Topaloglu. Robust assortment optimization under the multinomial
 logit choice model. *Operations Research*, 60(4):865–882, 2012.
- [46] A. Ruszczyński. Risk-averse dynamic programming for Markov decision processes. *Mathemat- ical Programming*, 125(2):235–261, 2010.
- [47] J. K. Satia and R. E. Lave Jr. Markovian decision processes with uncertain transition probabilities.
 Operations Research, 21(3):728–740, 1973.
- [48] A. Shapiro. Rectangular sets of probability measures. *Operations Research*, 64(2):528–541,
 2016.
- [49] A. Shapiro. Distributionally robust optimal control and MDP modeling. *Operations Research Letters*, 49(3):809–814, 2021.
- [50] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second
 edition, 2018.
- 477 [51] A. Tamar, Y. Glassner, and S. Mannor. Policy gradients beyond expectations: Conditional
 478 value-at-risk. Available on arXiv, 2014.
- [52] A. Tamar, S. Mannor, and H. Xu. Scaling up robust MDPs using function approximation. In
 Proceedings of the 31st International Conference of Machine Learning, 2014.
- [53] A. Tirinzoni, X. Chen, M. Petrik, and B. D. Ziebart. Policy-conditioned uncertainty sets for
 robust Markov decision processes. In *Advances in Neural Information Processing Systems*,
 volume 31, pages 8953–8963, 2018.
- [54] W. Wang and C. Lu. Projection onto the capped simplex. Available on arXiv, 2015.
- [55] Y. Wang and S. Zou. Online robust reinforcement learning with model uncertainty. In *Advances in Neural Information Processing Systems*, volume 34, pages 7193–7206, 2021.
- ⁴⁸⁷ [56] Yue Wang and Shaofeng Zou. Policy gradient method for robust reinforcement learning. In
 ⁴⁸⁸ Proceedings of the 39th International Conference on Machine Learning, pages 23484–23526,
 ⁴⁸⁹ 2022.
- 490 [57] T. Weissman, E. Ordentlich, G. Seroussi, S. Verdu, and M. J. Weinberger. In 491 equalities for the l₁ deviation of the empirical distribution. Technical Report,
 492 https://www.hpl.hp.com/research/info_theory/papers/HPL-2003-97R1Web.pdf,
 493 2003.
- 494 [58] W. Wiesemann, D. Kuhn, and B. Rustem. Robust Markov decision processes. *Mathematics of* 495 Operations Research, 38(1):153–183, 2013.
- 496 [59] H. Xiao, K. Yang, and X. Wang. Robust power control under channel uncertainty for cognitive
 radios with sensing delays. *IEEE Transactions on Wireless Communications*, 12(2):646–655,
 2013.

- [60] L. Xin and D. A. Goldberg. Distributionally robust inventory control when demand is a
 martingale. Available on arXiv, 2018.
- [61] H. Xu and S. Mannor. Distributionally robust Markov decision processes. In *Advances in Neural Information Processing Systems*, volume 23, pages 2505–2513, 2010.
- [62] H. Xu and S. Mannor. Distributionally robust Markov decision processes. *Mathematics of Operations Research*, 37(2):288–300, 2012.
- [63] P. Yu and H. Xu. Distributionally robust counterpart in Markov decision processes. *IEEE Transactions on Automatic Control*, 61(9):2538–2543, 2016.
- 507 [64] Y. Zhang, L. N. Steimle, and B. T. Denton. Robust Markov decision processes for medical
 508 treatment decisions. Available on Optimization Online, 2017.

509 Checklist

510	1.	For a	all authors
511		(a)	Do the main claims made in the abstract and introduction accurately reflect the paper's
512			contributions and scope? [Yes] The abstract explains that the main contributions of the
513			paper are (i) the development of a decomposition scheme that reduces the computation
514			of a seemingly unstructured robust Bellman operator to the repeated solution of highly
515			structured simplex projection problems and (ii) the fast solution of these simplex
516			projection problems for several classes of ϕ -divergences. These claims are backed up
517			in the introduction and the remainder of the paper.
518		(b)	Did you describe the limitations of your work? [Yes] We took great care to ensure that
519			our numerical results provide an objective and unbiased assessment of our solution
520			approach. In particular, we see that in one of the cases, the outperformance of our
521			approach over MOSEK slightly reduces for larger problem instances.
522		(c)	Did you discuss any potential negative societal impacts of your work? [N/A] Robust
523			MDPs are well-known in the literature, and our paper develops a new suite of fast
524			algorithms to solve these problems. As such, there are no new negative societal impacts
525			that we can identify.
526		(d)	Have you read the ethics review guidelines and ensured that your paper conforms to
527			them? [Yes] We have carefully read those guidelines, and to our best understanding
528			our paper fully complies with them.
529	2.	If yo	ou are including theoretical results
530		(a)	Did you state the full set of assumptions of all theoretical results? [Yes] All of our
531		()	results state the full set of assumptions, with exception of the blanket assumptions that
532			are assumed to hold throughout the paper and that are clearly marked as such.
533		(h)	Did you include complete proofs of all theoretical results? [Yes] All proofs are
534		(0)	contained in the appendix.
535	3.	If yo	bu ran experiments
536		(a)	Did you include the code, data, and instructions needed to reproduce the main experi-
537			mental results (either in the supplemental material or as a URL)? [Yes] All code, data
538			and instructions for our experimental results are published on GitHub. To maintain
539			anonymity during the review process, we do not provide a link in the current version of
540			the paper.
541		(b)	Did you specify all the training details (e.g., data splits, hyperparameters, how they
542			were chosen)? [Yes] All details are mentioned either in the numerical results section or
543			in the appendix.
544		(c)	Did you report error bars (e.g., with respect to the random seed after running experi-
545		. ,	ments multiple times)? [Yes] Included in the appendix.
546		(d)	Did you include the total amount of compute and the type of resources used (e.g., type
547			of GPUs, internal cluster, or cloud provider)? [Yes] Included in the numerical results
548			section.
549	4.	If yo	ou are using existing assets (e.g., code, data, models) or curating/releasing new assets
550		(a)	If your work uses existing assets did you cite the creators? [Yes] We use C ++ Python
551		(4)	and MOSEK, all of which are cited in the text.
552		(b)	Did you mention the license of the assets? [Yes] All licenses are mentioned.
553		(c)	Did vou include any new assets either in the supplemental material or as a URL? [Yes]
554		. /	All code, data and instructions for our experimental results are published on GitHub. To
555			maintain anonymity during the review process, we do not provide a link in the current
556			version of the paper.
557		(d)	Did you discuss whether and how consent was obtained from people whose data you're
558		. ,	using/curating? [N/A] We use synthetic data in our experiments.
559		(e)	Did you discuss whether the data you are using/curating contains personally identifiable
560		. /	information or offensive content? [N/A] We use synthetic data in our experiments.
561	5.	If yo	bu used crowdsourcing or conducted research with human subjects
	-	5.0	

(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] We use synthetic data in our experiments.
(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] We use synthetic data in our experiments.
(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] We use synthetic data in our experiments.