
Approximate Robust Control of Uncertain Dynamical Systems

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 This work studies the design of safe control policies for large-scale non-linear
2 systems operating in uncertain environments. In such a case, the robust control
3 framework is a principled approach to safety that aims to maximize the worst-case
4 performance of a system. However, the resulting optimization problem is generally
5 intractable for non-linear systems with continuous states. To overcome this issue,
6 we introduce two tractable methods that are based either on sampling or on a
7 conservative approximation of the robust objective. The proposed approaches are
8 applied to the problem of autonomous driving.

9 1 Introduction

10 Reinforcement Learning is a general framework that allows the optimal control of a Markov Decision
11 Process $(\mathcal{S}, \mathcal{A}, T, r)$ with state space \mathcal{S} , action space \mathcal{A} , reward function $r \in [0, 1]^{\mathcal{S} \times \mathcal{A}}$ and unknown
12 transition dynamics $T(s'|s, a) \in \mathcal{M}(\mathcal{S})^{\mathcal{S} \times \mathcal{A}}$ by searching for the policy $\pi \in \mathcal{M}(\mathcal{A})^{\mathcal{S}}$ with maximal
13 expected value v_{π}^T of the total discounted reward R_{π}^T :

$$R_{\pi}^T(s) \stackrel{\text{def}}{=} \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t), \quad v_{\pi}^T(s) \stackrel{\text{def}}{=} \mathbb{E}(R_{\pi}^T(s)), \quad (1)$$

14 where $s_0 = s$, $a_t \sim \pi(s_t)$, $s_{t+1} \sim T(s_{t+1}|s_t, a_t)$, $\gamma \in [0, 1)$ is the discount factor and $\mathcal{M}(X)$
15 denotes the set of probability measures over X .

16 Unfortunately, its application to real-world tasks has so far been limited by its considerable need for
17 experiences. It is generally recognized (Sutton, 1990; Atkeson and Santamaria, 1997) that the most
18 sample-efficient approach is the family of model-based methods which learn a nominal model \hat{T} of
19 the environment dynamics that is leveraged for policy search:

$$\max_{\pi} v_{\pi}^{\hat{T}} \quad (2)$$

20 One drawback of such methods is that they suffer from model bias; that is, they ignore the error
21 between the learned dynamics \hat{T} and the real environment T . It has been shown that model bias can
22 dramatically degrade the policy performances (Schneider, 1997).

23 Model errors can instead be explicitly considered and expressed through an *ambiguity set* of all
24 possible dynamics models. Such a set can be constructed from a history of observations by computing
25 the confidence regions associated with the system identification process (Iyengar, 2005; Nilim and El
26 Ghaoui, 2005; Dean et al., 2017; Maillard, 2017). In this work, we will consider ambiguity sets of
27 parametrized deterministic dynamical systems $s' = T_{\theta}(s, a)$ whose unknown parameters θ lie in a
28 compact set Θ of \mathbb{R}^p .

29 In the optimal control framework, model uncertainty is handled by maximizing the *expected* perfor-
30 mances with respect to unknown dynamics. In stark contrast, in real-world applications where failures
31 may turn out very costly, the decision maker often prefers to minimize the risk of the policy, which

32 can be defined with several metrics characterizing the distribution of the policy outcome (García and
 33 Fernández, 2015).

34 The robust control framework is a popular setting in which the risk of a policy is defined as the worst
 35 possible outcome realization among the ambiguity set, to guarantee a lower-bound performance of
 36 the robust policy when executed on the true model:

$$\max_{\pi} \min_T v_{\pi}^T \quad (3)$$

37 Robust optimization has been studied in the context of finite Markov Decision Processes (MDP) with
 38 uncertain parameters by Iyengar (2005), Nilim and El Ghaoui (2005) and Wiesemann et al. (2013).
 39 They show that the main results of Dynamic Programming can be extended to their robust counterparts
 40 only when the dynamics ambiguity set verifies certain rectangularity properties. In the control theory
 41 community, the robust control problem is mainly restricted to the context of linear dynamical systems
 42 with bounded uncertainty in the time or frequency domain, where the objective is to guarantee
 43 stability (e.g. \mathcal{H}_{∞} -optimal control, see Basar and Bernhard, 1996) or performance (e.g. LQ optimal
 44 control theory, see Petersen and Tempo, 2014). The existing nonlinear robust control approaches
 45 such as sliding mode control (Li et al., 2017), feedback linearization, backstepping, passivation and
 46 input-to-states stabilization (Khalil, 2014) are usually based on canonical representations of regulated
 47 dynamics and admit constructive numeric realizations for systems of rather low dimensions.

48 There have been few attempts of robust control of large-scale systems with both continuous states
 49 and non-linear dynamics, which is the focus of this paper. Our contribution is twofold. In section 2,
 50 we first consider a simpler case where the ambiguity set Θ and action space \mathcal{A} are both finite and
 51 introduce a sampling-based planner that approximately maximizes the robust objective (3). In section
 52 3, we move to continuous ambiguity sets and form a conservative relaxation of the robust policy
 53 evaluation problem using interval predictors. In section 4, we illustrate the benefits of both techniques
 54 (for discrete, versus continuous Θ) on a problem of tactical decision-making for autonomous driving.

55 2 Sampling-based planning

56 If the true dynamics model T_{θ} were known and the action-space \mathcal{A} finite, sampling-based algorithms
 57 could be used to perform approximate optimal planning. In order to generalize to the robust setting,
 58 we need to make the following assumption about the structure of the ambiguity set:

59 **Assumption 1** (Structure). *The ambiguity set Θ and the action space \mathcal{A} are discrete and finite:*

$$\mathcal{A} = \{a_k\}_{k \in [1, K]} \quad \text{and} \quad \Theta = \{\theta_m\}_{m \in [1, M]} \quad (4)$$

60 We slightly abuse notation and denote $T_m = T_{\theta_m}$.

61 Such a structure of the ambiguity set typically stems directly from expert knowledge of the problem
 62 at hand. In general, it is nonrectangular, which implies that the Robust Bellman Equation does not
 63 hold (Wiesemann et al., 2013). This prevents us from building on planners that implicitly use this
 64 property and generate trajectories step-by-step by picking promising successor states, such as MCTS
 65 (Coulom, 2006) or UCT (Kocsis and Szepesvári, 2006). Instead, we turn to algorithms that perform
 66 optimistic sampling of entire sequences of actions and work directly at the leaves of the expanded
 67 tree (see, e.g. Bubeck and Munos, 2010). More precisely, we build on the work of Hren and Munos
 68 (2008) on optimistic planning for deterministic dynamics, which we extend to the robust setting.

69 We use similar notations and consider the infinite look-ahead tree \mathcal{T} composed of all reachable states.
 70 Each node corresponds to a joint state $\{s_{m,t}\}_{m \in [1, M]}$ associated with the different dynamics T_m .
 71 The root starts at the current state, and all nodes have K children, each corresponding to an action
 72 $a_k \in \mathcal{A}$ and associated with the successor joint state $\{s_{m,t+1} = T_m(s_{m,t}, a_k)\}_{m \in [1, M]}$. We use the
 73 standard notations over alphabets to refer to nodes in \mathcal{T} as action sequences. Thus, a finite word
 74 $i \in \mathcal{A}^*$ of length d represents the node obtained following the action sequence (i_0, \dots, i_d) from
 75 the root. Sequences $i \in \mathcal{A}^*$ and $j \in \mathcal{A}^*$ can be concatenated as $ij \in \mathcal{A}^*$, the set of suffixes of i is
 76 $i\mathcal{A}^{\infty} = \{j \in \mathcal{A}^{\infty} : \exists h \in \mathcal{A}^{\infty} \text{ such that } j = ih\}$, and the empty sequence is \emptyset .

77 The sample complexity is expressed in terms of number n of expanded nodes. It is related to the
 78 number of calls to dynamics models: when a node i is expanded, all successor states are computed
 79 for all K actions and M dynamics. At an iteration n , we denote \mathcal{T}_n the tree of already expanded
 80 nodes, and \mathcal{L}_n the set of its leaves.

81 **Definition** Fix a dynamics model $m \in [1, M]$. Hren and Munos (2008) define for any node $i \in \mathcal{T}$
 82 of depth d the optimal value v_i^m , its lower bound u-value u_i^m and upper-bound b-value b_i^m . These
 83 variables depend on the dynamics m and will therefore be referred to with a superscript m notation.

84 We extend these dynamics-dependent variables to the robust setting, using superscript r in notations.

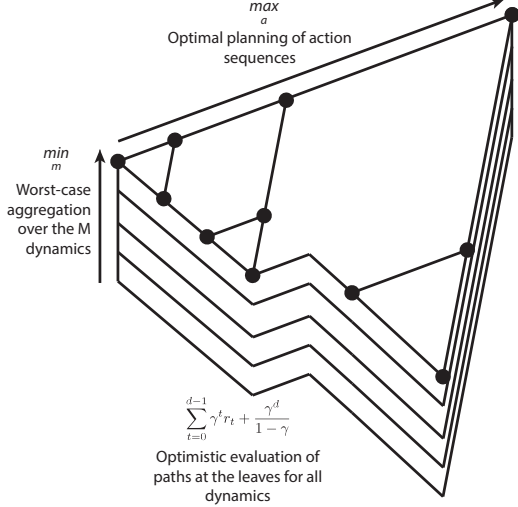


Figure 1: The computation of robust b-values in Algorithm 1. The simulation of trajectories for every dynamics model T_m is represented as stacked versions of the expanded tree \mathcal{T}_n .

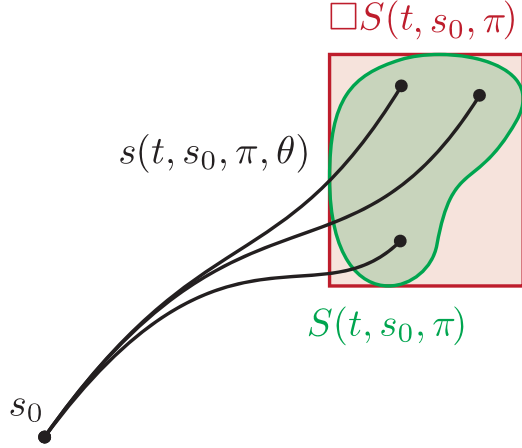


Figure 2: A few trajectories are sampled from an initial state s_0 following a policy π with various dynamics parameters θ_m (in black). The union of reachability sets is shown in green, and its interval hull in red.

- 85 • The robust value v_i^r of a path $i \in \mathcal{A}^*$ as the restriction of (3) to policies that start with the
86 action sequence i :
$$v_i^r \stackrel{\text{def}}{=} \max_{\pi \in i\mathcal{A}^\infty} \min_{m \in [1, M]} R_\pi^{T_m} \quad (5)$$

87 By definition, the robust value of (3) is recovered at the root $v_\emptyset^r = v^r$.
88 Moreover, for $i \in \mathcal{T}_n \setminus \mathcal{L}_n$ we have

$$v_i^r = \max_{\pi \in i\mathcal{A}^\infty} \min_{m \in [1, M]} R_\pi^{T_m} = \max_{a \in \mathcal{A}} \max_{\pi \in ia\mathcal{A}^\infty} \min_{m \in [1, M]} R_\pi^{T_m} = \max_{a \in \mathcal{A}} v_{ia}^r \quad (6)$$

- 89 • The robust u-value u_i^r of a leaf node $i \in \mathcal{L}_n$ is the worst-case discounted sum of rewards
90 $r_t = r(s_{m,t}, i_t)$ from the root to i . It is then backed-up to the rest of the tree:

$$u_i^r(n) \stackrel{\text{def}}{=} \begin{cases} \min_{m \in [1, M]} \sum_{t=0}^{d-1} \gamma^t r_t & \text{if } i \in \mathcal{L}_n ; \\ \max_{a \in \mathcal{A}} u_{ia}^r(n) & \text{if } i \in \mathcal{T}_n \setminus \mathcal{L}_n \end{cases} \quad (7)$$

- 91 • Likewise, the robust b-value b_i^r is defined at leaf nodes and backed-up to the rest of the tree:

$$b_i^r(n) \stackrel{\text{def}}{=} \begin{cases} u_i^r(n) + \frac{\gamma^d}{1-\gamma} & \text{if } i \in \mathcal{L}_n ; \\ \max_{a \in \mathcal{A}} b_{ia}^r(n) & \text{if } i \in \mathcal{T}_n \setminus \mathcal{L}_n \end{cases} \quad (8)$$

92 An illustration of the computation of the robust b-values is presented in Figure 1.

93 **Remark 1** (On the ordering of min and max). *In the definition of $u_i^r(n)$ it is essential that the*
94 *minimum among the models is only taken at the end of trajectories, in the same way as for the robust*
95 *objective (3) in which the worst-case dynamics is only determined after the policy has been fully*
96 *specified. Assume that $u_i^r(n)$ is instead naively defined as:*

$$u_i^r(n) = \min_{m \in [1, M]} u_i^m(n),$$

This would not recover the robust policy, as we show in Figure 3 with a simple counter-example.

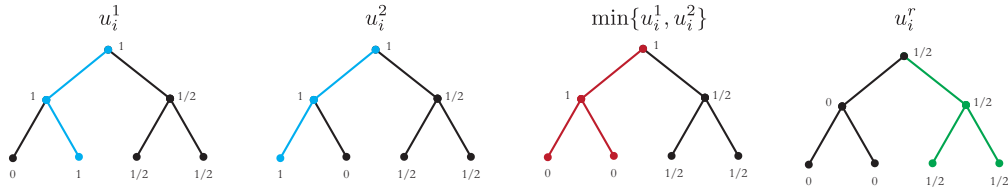


Figure 3: From left to right: two simple models and corresponding u-values with optimal sequences in blue; the naive version of the robust values returns sub-optimal paths in red; our robust u-value properly recovers the robust policy in green.

97

Algorithm 1: Deterministic Robust Optimistic Planning

```

1 Initialize  $\mathcal{T}$  to a root and expand it. Set  $n = 1$ .
2 while Numerical resource available do
3   Compute the robust u-values  $u_i^r(n)$  and robust b-values  $b_i^r(n)$ .
4   Expand  $\operatorname{argmax}_{i \in \mathcal{L}_n} b_i^r(n)$ .
5    $n = n + 1$ 
6 return  $\operatorname{argmax}_{a \in \mathcal{A}} u_a^r(n)$ 

```

98 From these definitions we introduce Algorithm 1, and analyse its sample-efficiency in Theorem 1.

99 **Lemma 1** (Robust values ordering). *The robust values, u-values and b-values exhibit similar properties as the optimal values, u-values and b-values, that is: for all $0 < t < n$ and $i \in \mathcal{T}_n$,*

$$u_i^r(t) \leq u_i^r(n) \leq v_i^r \leq b_i^r(n) \leq b_i^r(t) \quad (9)$$

101

102 *Proof.* This result stems directly from the definitions, see more details in Appendix A.1. ■

103 The simple regret of the action a returned by Algorithm 1 after n rounds is defined as:

$$\mathcal{R}_n = v^r - v_a^r \quad (10)$$

104 We will say that $\mathcal{R}_n = O(\varepsilon)$ for some $\varepsilon > 0$ if there exist $\rho > 0$ and $n_0 > 0$ such that $\mathcal{R}_n \leq \rho\varepsilon$ for
 105 all $n \geq n_0$. A node $i \in \mathcal{T}$ is said to be ε -optimal, in a robust sense, if and only if $v_i^r \geq v^r - \varepsilon$ for
 106 some $\varepsilon > 0$. The proportion of ε -optimal nodes at depth d is then defined as $p_d(\varepsilon) = |\{i \in \mathcal{A}^d \text{ s.t. } i \text{ is}$
 107 ε -optimal\}| / K^d . Further we will assume that for the graph \mathcal{T} the following hypothesis is satisfied:

108 **Assumption 2** (Proportion of near-optimal nodes). *There exist $\beta \in [0, \frac{\log K}{\log 1/\gamma}]$, $c > 0$ and $d_0 > 0$
 109 such that $p_d(\varepsilon) \leq c\varepsilon^\beta$ for all $\varepsilon > 0$ and $d \geq d_0$.*

110 **Theorem 1** (Regret bound). *Let $\kappa = K\gamma^\beta \in [1, K]$. Then the simple regret of Algorithm 1 is:*

$$\text{If } \kappa > 1, \quad \mathcal{R}_n = O\left(n^{-\frac{\log 1/\gamma}{\log \kappa}}\right) \quad (11)$$

$$\text{If } \kappa = 1, \quad \mathcal{R}_n = O\left(\gamma^{\frac{(1-\gamma)^\beta}{c} n}\right) \quad (12)$$

111 *Proof.* We use the properties shown in Lemma 1 and derive a robust counterpart of the proof of Hren
 112 and Munos (2008), which we only modify slightly. See more details in Appendix A.2 ■

113 3 Interval predictors

114 In this section, we assume that the ambiguity set Θ is continuous and bounded.

115 In the robust objective (3), the min operator only requires us to describe the set of states that can be
 116 reached with non-zero probability.

117 **Definition** The **reachability set** S at time t is the set of all states that can be reached by starting
 118 from initial state $s_0 \in \mathcal{S}$ and following a policy $\pi \in \mathcal{A}^S$ along the transition dynamics $T_\theta \in \mathcal{S}^{\mathcal{S} \times \mathcal{A}}$.

$$S(t, s_0, \pi) \stackrel{\text{def}}{=} \{s_t : \exists \theta \in \Theta \text{ s.t. } s_{k+1} = T_\theta(s_k, a_k), a_k = \pi(s_k), k = 0, \dots, t-1\} \quad (13)$$

120 This set can still have a complex shape. We approximate it by an overset easier to represent and
 121 manipulate: its interval hull.

122 **Definition** The **interval hull** of S , denoted $\square S = [\underline{s}, \bar{s}]$ is the smallest interval containing it:

$$\underline{s}(t, s_0, \pi) \stackrel{\text{def}}{=} \min S(t, s_0, \pi) \quad \bar{s}(t, s_0, \pi) \stackrel{\text{def}}{=} \max S(t, s_0, \pi) \quad (14)$$

123 The max and min operators are applied element-wise. This set is illustrated in Figure 2.

124 State intervals $\square S$ have been used to describe the evolution of uncertain systems and derive feedback
 125 laws that achieve closed-loop stability in the presence of bounded disturbances (Stinga and Bunciu,
 126 2012; Efimov and Raïssi, 2016; Dinh and Ito, 2017).

127 The main techniques of interval simulation have been listed and described in a survey by Puig et al.
 128 (2005), in which they are sorted into two categories. Region-based methods use the estimate of $\square S$

Algorithm 2: Interval-based Robust Control

```
1 Algorithm robust_control( $s_0$ )
2   Initialize a set  $\Pi$  of policies
3   while resources available do
4     evaluate() each policy  $\pi \in \Pi$  at current state  $s_0$ 
5     Update  $\Pi$  by policy search
6   end
7   return  $\operatorname{argmax}_{\pi \in \Pi} \hat{v}^r(\pi)$ 
1 Procedure evaluate( $\pi, s_0$ )
2   Compute the state interval  $\square S(t, s_0, \pi)$  on a horizon  $t \in [0, H]$ 
3   Minimize  $r$  over the intervals  $\square S(t, s_0, \pi)$  for all  $t \in [0, H]$ 
4   return  $\hat{v}^r(\pi)$ 
```

129 at previous timestep $t - 1$ to bootstrap the current estimate at time t . They are based on application
130 of the theory of positive systems, which are frequently computationally efficient. However, the
131 positive inclusion dynamics of a system may lead to overestimations of the true $\square S$ and even unstable
132 behaviour. Trajectory-based methods estimate $\square S$ by taking the max and min in (14) over sampled
133 trajectories for $\theta \in \Theta$. These methods produce subset estimates of the true $\square S$, do not suffer from
134 the wrapping effect, but are often more computationally costly.

135 In this work, we leverage them to derive a proxy for the robust objective (3).

136 **Definition** Let us denote the robust objective of equation (3) as $v^r(\pi) \stackrel{\text{def}}{=} \min_{\theta \in \Theta} v_{\pi}^{T_{\theta}}$.

137 We define the **surrogate objective** \hat{v}^r on a finite horizon $H > 0$ as:

$$\hat{v}^r(\pi) \stackrel{\text{def}}{=} \sum_{t=0}^H \gamma^t \min_{s \in \square S(t, s_0, \pi)} r(s, \pi(s)) \quad (15)$$

138 **Property 1** (Lower bound). *The surrogate objective \hat{v}^r is a lower bound of the true objective v^r :*

$$\forall \pi, \hat{v}^r(\pi) \leq v^r(\pi) \quad (16)$$

140 *Proof.* By bounding the collected rewards by their minimum over $\square S(t)$. See Appendix A.3 ■

141 The robust objective error $v^r - \hat{v}^r$ stems from two terms: the interval approximation of the reachable
142 set and the loss of time-dependency between the states within a single trajectory. If both these
143 approximations are tight enough, maximizing the lower bound \hat{v}^r will increase the true objective
144 v^r , which is the idea behind Algorithm 2. It is classically structured as an alternation of a Policy
145 Evaluation step, during which the surrogate objective $\hat{v}^r(\pi)$ is evaluated for a set of policies Π , and
146 a Policy Search step which aims to steer the set of policies Π towards regions where the surrogate
147 objective is maximal. The main Policy Search algorithms are listed in a survey by Deisenroth
148 (2011). In this case, derivative-free methods such as evolutionary strategies (e.g. CMAES) would be
149 more appropriate than policy gradient methods, since \hat{v}^r cannot be easily differentiated. Planning
150 algorithms can also be used to exploit the dynamics and structure of the surrogate objective.

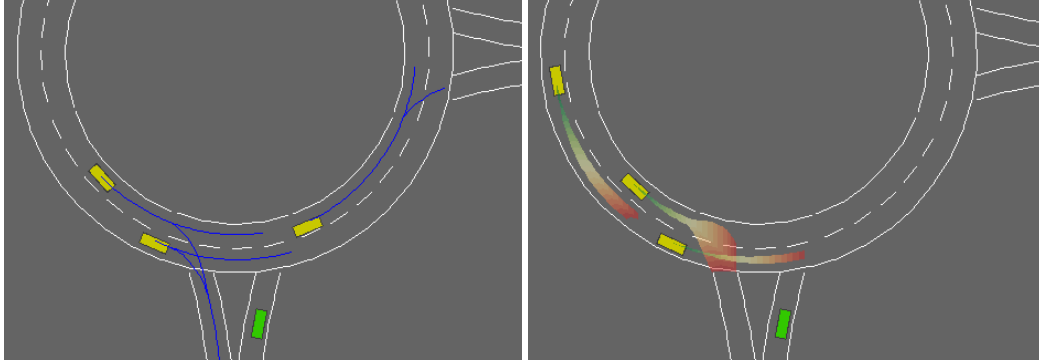
151 4 Experiments

152 Most autonomous driving architectures perform sequentially the prediction of other drivers' trajec-
153 tories and the planning of a collision-free path for the ego-vehicle. As a consequence, they fail
154 to account for interactions between the traffic participants and the ego-vehicle, leading to overly
155 conservative decisions and a lack of negotiation abilities (Trautman and Krause, 2010). In this
156 work, we perform both tasks *jointly* to anticipate the effect of our own decisions on the dynamics
157 of the nearby traffic. But human decisions are not fully predictable and cannot be reduced to a
158 single deterministic model. To avoid model bias, we provide a whole ambiguity set of reasonable
159 closed-loop behavioural models for other vehicles, and plan robustly with respect to this ambiguity.

160 We introduce a new environment for simulated highway driving and tactical decision-making.¹

161 Vehicle motion is described by the Kinematic Bicycle Model (see, e.g. Polack et al., 2017). They
162 follow a lane keeping lateral behaviour, and a longitudinal behaviour inspired by the Intelligent Driver
163 Model (Treiber et al., 2000) which balances reaching a desired velocity and respecting a safe time

¹Source code is available at <https://github.com/eleurent/highway-env>



(a) The possible trajectories (blue) for fixed behaviours and varying destinations (b) The possible trajectories (green-red gradient) for fixed destination and varying behaviours

Figure 4: The highway-env environment. The ego-vehicle (green) is approaching a roundabout with flowing traffic (yellow).

Table 1: Performances of robust planners on two ambiguous environments.

Ambiguity set	Agent	Worst-case return	Mean return \pm std
True model	Oracle	9.83	10.84 ± 0.16
Discrete	Nominal	2.09	8.85 ± 3.53
	Algorithm 1	8.99	10.78 ± 0.34
Continuous	Nominal	1.99	9.95 ± 2.38
	Algorithm 2	7.88	10.73 ± 0.61

164 gap. The lane-change decisions are determined by the MOBIL model (Kesting et al., 2007): they
 165 must increase the vehicles accelerations while satisfying safe braking decelerations. The behaviour
 166 parameters θ of each traffic participant are sampled uniformly from a set Θ .

167 The ego-vehicle can be controlled with a finite set of tactical decisions $\mathcal{A} = \{\text{no-op, right-lane,}$
 168 $\text{left-lane, faster, slower}\}$ implemented by low-level controllers. It is rewarded for driving fast
 169 along a planned route while avoiding collisions. More information on the environment modelling is
 170 provided in the appendices.

171 We carry out two experiments²: First, the behavioural parameters of traffic participants are fixed but
 172 their planned routes are unknown: we enumerate every direction they can take at their next intersection
 173 (see Figure 4a) and plan robustly with respect to this finite ambiguity set using Algorithm 1. Second,
 174 we assume on the contrary that the agents’ planned routes are known but not their behavioral
 175 parameters (see Figure 4b). We plan robustly with respect to this continuous ambiguity set using
 176 Algorithm 2. Crucially, the state intervals prediction is conditioned on the planned policy π .

177 In both experiments, we compare the performance of the robust planner to an oracle model that has
 178 perfect knowledge of the systems dynamics, and to a nominal planner that plans optimistically with
 179 respect to a dynamics model sampled uniformly from the ambiguity set. Statistics are collected from
 180 100 episodes with random environment initialization. Results are presented in Table 1.

181 5 Conclusion

182 This paper has presented two methods for approximately solving the robust control problem. In
 183 the simpler case of finite ambiguity set and action space, we use optimistic planning and provide
 184 an upper bound for the simple regret. A direct consequence is that we recover the robust policy as
 185 the computational budget increases. In the general case, we use interval prediction to efficiently
 186 solve a conservative approximation of the robust objective while providing a lower bound for the
 187 performance of a policy when applied to the unknown true model. However, this method is lossy and
 188 does not enjoy asymptotic consistency. Both algorithms are flexible, allowing to handle a variety of
 189 parametrized dynamical systems, and practical, with a focus on computational efficiency. The two
 190 methods are also orthogonal, which means they can be combined to deal with complex ambiguity
 191 sets that display both continuous and discrete features, such as disjoint unions of connected sets.

²Video and source code are available at <https://eleurent.github.io/robust-control/>

References

- 192
193 C.G. Atkeson and J.C. Santamaria. A comparison of direct and model-based reinforcement learning. *Proceedings*
194 *of International Conference on Robotics and Automation*, 1997.
- 195 T. Basar and P. Bernhard. *H infinity - Optimal Control and Related Minimax Design Problems: A Dynamic*
196 *Game Approach*, volume 41. 1996.
- 197 Sébastien Bubeck and Rémi Munos. Open Loop Optimistic Planning. Technical report, 2010.
- 198 Rémi Coulom. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. pages 72–83, 2006.
- 199 Sarah Dean, Horia Mania, Nikolai Matni, Benjamin Recht, and Stephen Tu. On the Sample Complexity of the
200 Linear Quadratic Regulator. 2017.
- 201 Marc Peter Deisenroth. A Survey on Policy Search for Robotics. *Foundations and Trends in Robotics*, 2011.
- 202 Thach Ngoc Dinh and Hiroshi Ito. Decentralization of Interval Observers for Robust Controlling and Monitoring
203 a Class of Nonlinear Systems. 10:117–123, 2017.
- 204 D. Efimov and T. Raïssi. Design of interval observers for uncertain dynamical systems. *Automation and Remote*
205 *Control*, 77:191–225, 2016.
- 206 Javier García and Fernando Fernández. A Comprehensive Survey on Safe Reinforcement Learning. *Journal of*
207 *Machine Learning Research*, 16:1437–1480, 2015. ISSN 15337928.
- 208 Jean-François Hren and Rémi Munos. Optimistic planning for deterministic systems. 2008.
- 209 Garud N. Iyengar. Robust Dynamic Programming. *Mathematics of Operations Research*, 30:257–280, 2005.
- 210 Arne Kesting, Martin Treiber, and Dirk Helbing. General Lane-Changing Model MOBIL for Car-Following
211 Models. *Transportation Research Record: Journal of the Transportation Research Board*, pages 86–94, 2007.
- 212 Hassan K. Khalil. *Nonlinear Control*. Pearson, 2014. ISBN 013349926X.
- 213 Levente Kocsis and Csaba Szepesvári. Bandit Based Monte-Carlo Planning. pages 282–293, 2006.
- 214 Shihua Li, Xinghuo Yu, Leonid Fridman, Zhihong Man, and Xiangyu Wang. *Advances in Variable Structure*
215 *Systems and Sliding Mode Control – Theory and Applications (Studies in Systems, Decision and Control)*.
216 Springer, 2017.
- 217 Odalric-Ambrym Maillard. Self-normalization techniques for streaming confident regression. 2017.
- 218 Arnab Nilim and Laurent El Ghaoui. Robust Control of Markov Decision Processes with Uncertain Transition
219 Matrices. *Operations Research*, 53:780–798, 2005.
- 220 Ian R. Petersen and Roberto Tempo. Robust control of uncertain systems: Classical results and recent develop-
221 ments. *Automatica*, 50(5):1315–1335, 2014.
- 222 Philip Polack, Florent Alché, and Brigitte D’Andréa-Novel. The Kinematic Bicycle Model : a Consistent Model
223 for Planning Feasible Trajectories for Autonomous Vehicles ? pages 6–8, 2017.
- 224 Vicenç Puig, Alexandru Stancu, and Joseba Quevedo. *Simulation of Uncertain Dynamic Systems Described By*
225 *Interval Models: a Survey*, volume 38. IFAC, 2005.
- 226 JG Schneider. Exploiting model uncertainty estimates for safe dynamic control learning. *Advances in neural*
227 *information processing systems*, pages 1047—1053, 1997.
- 228 F. Stinga and E. Bunciu. Robust interval observer and nonlinear predictive control of an active sludge process.
229 *System Theory, Control and Computing*, (1), 2012.
- 230 Richard S. Sutton. Integrated Architectures for Learning, Planning, and Reacting Based on Approximating
231 Dynamic Programming. In *Machine Learning Proceedings 1990*. 1990.
- 232 Peter Trautman and Andreas Krause. Unfreezing the robot: Navigation in dense, interacting crowds. In *IEEE/RSJ*
233 *2010 International Conference on Intelligent Robots and Systems, IROS 2010 - Conference Proceedings*,
234 2010.
- 235 Martin Treiber, Ansgar Hennecke, and Dirk Helbing. Congested traffic states in empirical observations and mi-
236 croscopic simulations. *Physical Review E - Statistical Physics, Plasmas, Fluids, and Related Interdisciplinary*
237 *Topics*, 62(2):1805–1824, 2000.
- 238 Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov Decision Processes. . . . *of Operations*
239 *Research*, pages 1–52, 2013.

Supplementary material

240

241 A Detailed proofs

242 A.1 Lemma 1

243 *Proof.* By definition, when starting with sequence i , the value $u_i^m(n)$ represents the minimum admissible
 244 reward, while $b_i^m(n)$ corresponds to the best admissible reward achievable with respect to the the possible
 245 continuations of i . Thus, for all $i \in \mathcal{A}^*$, $u_i^m(n)$ and $u_i^r(n)$ are non-decreasing functions of n and $b_i^m(n)$ and
 246 $b_i^r(n)$ are a non-increasing functions of n , while v_i^m and v_i^r do not depend on n .

247 Moreover, since the reward function r is assumed to have values in $[0, 1]$, the sum of discounted rewards from a
 248 node of depth d is at most $\gamma^d + \gamma^{d+1} + \dots = \frac{\gamma^d}{1-\gamma}$. As a consequence, for all $n \geq 0$, $i \in \mathcal{L}_n$ of depth d , and
 249 any sequence of rewards $(r_t)_{t \in \mathbb{N}}$ obtained from following a path in $i\mathcal{A}^\infty$ with any dynamics $m \in [1, M]$:

$$\sum_{t=0}^{d-1} \gamma^t r_t \leq \sum_{t=0}^{d-1} \gamma^t r_t + \sum_{t=d}^{\infty} \gamma^t r_t \leq \sum_{t=0}^{d-1} \gamma^t r_t + \frac{\gamma^d}{1-\gamma}$$

250 That is equivalent to:

$$u_i^m(n) \leq \sum_{t=0}^{\infty} \gamma^t r_t \leq b_i^m(n)$$

251 Hence,

$$\min_{m \in [1, M]} u_i^m(n) \leq \min_{m \in [1, M]} \sum_{t=0}^{\infty} \gamma^t r_t \leq \min_{m \in [1, M]} b_i^m(n) \quad (17)$$

252 And as the left-hand and right-hand sides of (17) are independent of the particular path that was followed in
 253 $i\mathcal{A}^\infty$, it also holds for the robust path:

$$\min_{m \in [1, M]} u_i^m(n) \leq \max_{\pi \in i\mathcal{A}^\infty} \min_{m \in [1, M]} \sum_{t=0}^{\infty} \gamma^t r_t \leq \min_{m \in [1, M]} b_i^m(n)$$

254 that is,

$$u_i^r(n) \leq v_i^r \leq b_i^r(n) \quad (18)$$

255 Finally, (18) is extended to the rest of \mathcal{T}_n by recursive application of (6), (7) and (8). \blacksquare

256 A.2 Theorem 1

257 *Proof.* Hren and Munos (2008) first show in Theorem 2 that the simple regret of their optimistic planner is
 258 bounded by $\frac{\gamma^{d_n}}{1-\gamma}$ where d_n is the depth of \mathcal{T}_n . This properties relies on the fact that the returned action belongs
 259 to the deepest explored branch, which we can show likewise by contradiction using Lemma 1. This yields
 260 directly that $a = i_0$ where i is some node of maximal depth d_n expanded at round $t \leq n$, which by Algorithm 1
 261 verifies $b_a^r(t) = b_i^r(t) = \max_{x \in \mathcal{A}} b_x^r(t)$ and:

$$v^r - v_a^r = v_{a^*}^r - v_a^r \leq b_{a^*}^r(t) - v_a^r \leq b_a^r(t) - u_a^r(t) = b_i^r(t) - u_i^r(t) = \frac{\gamma^{d_n}}{1-\gamma} \quad (19)$$

262 Secondly, they bound the depth d_n of \mathcal{T}_n with respect to n . To that end, they show that the expanded nodes
 263 always belong to the sub-tree \mathcal{T}_∞ of all the nodes of depth d that are $\frac{\gamma^d}{1-\gamma}$ -optimal. Indeed, if a node i of
 264 depth d is expanded at round n , then $b_i^r(n) \geq b_j^r(n)$ for all $j \in \mathcal{L}_n$ by Algorithm 1, thus the max-backups
 265 of (8) up to the root yield $b_i^r(n) = b_\emptyset^r(n)$. Moreover, by Lemma 1 we have that $b_\emptyset^r(n) \geq v_\emptyset^r = v^r$ and so
 266 $v_i^r \geq u_i^r(n) = b_i^r(n) - \frac{\gamma^d}{1-\gamma} \geq v^r - \frac{\gamma^d}{1-\gamma}$, thus $i \in \mathcal{T}_\infty$.

267 Then from Assumption 2 and the definition of β applied to nodes in \mathcal{T}_∞ , there exists d_0 and c such that the
 268 number n_d of nodes of depth $d \geq d_0$ in \mathcal{T}_∞ is bounded by $c \left(\frac{\gamma^d}{1-\gamma} \right)^\beta K^d$. As a consequence,

$$\begin{aligned} n &= \sum_{d=0}^{d_n} n_d = n_0 + \sum_{d=d_0+1}^{d_n} n_d \\ &\leq n_0 + \sum_{d=d_0+1}^{d_n} c \left(\frac{\gamma^d}{1-\gamma} \right)^\beta K^d \\ &= n_0 + c' \sum_{d=d_0+1}^{d_n} \kappa^d \end{aligned}$$

269 where $c' = \frac{c}{(1-\gamma)^\beta}$.

- 270 • If $\kappa > 1$, then $n \leq n_0 + c' \kappa^{d_0+1} \frac{\kappa^{d_n-d_0}-1}{\kappa-1}$ and thus $d_n \geq d_0 + \log_{\kappa} \frac{(n-n_0)(\kappa-1)}{c' \kappa^{d_0+1}}$. We conclude
 271 from (19) that $\mathcal{R}_n \leq \frac{\gamma^{d_n}}{1-\gamma} = \frac{1}{1-\gamma} \left(\frac{(n-n_0)(\kappa-1)}{c' \kappa^{d_0+1}} \right)^{\frac{\log \gamma}{\log \kappa}} = O \left(n^{-\frac{\log 1/\gamma}{\log \kappa}} \right)$.
- 272 • If $\kappa = 1$, then $n \leq n_0 + c'(d_n - d_0)$, hence from (19) we have $\mathcal{R}_n = O(\gamma^{nc'})$.

273

274 A.3 Property 1

275 *Proof.* For any $\theta \in \Theta$, $t \in [0, H]$ and any trajectory (s_0, \dots, s_t) sampled from π and T_{θ} ,

$$s_t \in S(t, s_0, \pi) \subset \square S(t, s_0, \pi)$$

276 Hence,

$$R_{\pi}^{T_{\theta}} = \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \geq \sum_{t=0}^H \gamma^t r(s_t, a_t) \geq \sum_{t=0}^H \min_{s \in \square S(t, s_0, \pi)} \gamma^t r(s, \pi(s)) = \hat{v}^r(\pi)$$

277 And finally,

$$v^r(\pi) = \min_{\theta \in \Theta} v_{\pi}^{T_{\theta}} = \min_{\theta \in \Theta} \mathbb{E}(R_{\pi}^{T_{\theta}}) \geq \hat{v}^r(\pi)$$

278

279 B Environment dynamics

280 B.1 Kinematics

281 The vehicles kinematics are represented by the Kinematic Bicycle Model:

$$\dot{x} = v \cos(\psi), \quad (20)$$

$$\dot{y} = v \sin(\psi), \quad (21)$$

$$\dot{v} = a, \quad (22)$$

$$\dot{\psi} = \frac{v}{l} \tan(\beta), \quad (23)$$

282 where (x, y) is the vehicle position, v its forward velocity and ψ its heading, l is the vehicle half-length, a is the
 283 acceleration command and β is the slip angle at the center of gravity, used as a steering command.

284 Each vehicle i is represented by its kinematics $X_i = [x_i, y_i, v_i, \psi_i]$. The joint state is represented by $s =$
 285 $\{X_1, \dots, X_N\}$

286 B.2 Longitudinal control

287 The acceleration control is assumed to be linearly parametrized:

$$a = \theta_a^T \phi_a(s, i), \quad (24)$$

288 where θ_a is an uncertain weight vector, and $\phi_a(s, i)$ is a feature vector that depends on the joint state s and
 289 considered vehicle i .

290 It is composed of:

- 291 • a target velocity seeking term,
 292 • a braking term to adjust velocity w.r.t. the front vehicle ,
 293 • a braking term to respect a safe distance w.r.t. the front vehicle.

294 Denoting f_i the front vehicle preceding vehicle i , ϕ_a is defined by

$$\phi_a(s, i) = \begin{bmatrix} v_0 - v_i \\ n(v_{f_i} - v_i) \\ n(x_{f_i} - x_i - (d_0 + v_i T)) \end{bmatrix} \quad (25)$$

295 where n is the negative part function $n(x) = \min(x, 0)$ and v_0 , d_0 and T respectively denote the speed limit,
296 jam distance and time gap given by traffic rules.

297 We observe that this model exhibits similar qualitative behaviours to the IDM's.

298 **B.3 Lateral control**

299 A non-linear lane-keeping controller is implemented as follows: a lane L with lateral position y_L and heading
300 ψ_L is tracked by performing

301 1. Position control

$$v_{y_{cmd}} = K_{p_y}(y_L - y) \quad (26)$$

302 2. Lateral velocity to heading conversion

$$\psi_{ref} = \psi_L + \sin^{-1}\left(\frac{v_{y_{cmd}}}{v}\right) \quad (27)$$

303 3. Heading control

$$\psi_{cmd} = K_{p_\psi}(\psi_{ref} - \psi) \quad (28)$$

304 4. Heading rate to steering angle conversion

$$\beta = \tan^{-1}\left(\frac{l}{v}\psi_{cmd}\right) \quad (29)$$

305 Finally,

$$\beta = \tan^{-1}\left(\frac{l}{v}K_{p_\psi}\left(\psi_L + \sin^{-1}\left(K_{p_y}\frac{y_L - y}{v}\right) - \psi\right)\right) \quad (30)$$

306 This non-linear controller presented in subsection can be linearised around its equilibrium $(y, \psi) = (y_L, \psi_L)$.

$$\frac{l}{v} \tan \beta = K_{p_\psi}\left(\psi_L + \sin^{-1}\left(K_{p_y}\frac{y_L - y}{v}\right) - \psi\right) \quad (31)$$

$$\simeq \frac{l}{v}\left(K_{p_\psi}\left(\psi_L + \left(K_{p_y}\frac{y_L - y}{v}\right) - \psi\right)\right) \quad (32)$$

$$= \theta_b^T \phi_b \quad (33)$$

307 with

$$\theta_b = [K_{p_\psi} \quad K_{p_y}K_{p_\psi}]^T \quad (34)$$

308 and

$$\phi_b = \begin{bmatrix} \psi_L - \psi \\ \frac{l}{v}(y_L - y) \end{bmatrix} \quad (35)$$

309 **B.4 Discrete behaviour**

310 The MOBIL model (Kesting et al., 2007), which stands for *Minimizing Overall Braking Induced by Lane*
311 *Changes*, is a discrete lateral decision model that formulates a criterion for lane changes in terms of safe braking
312 decelerations and increased overall accelerations according to a longitudinal model.

313 It states that a lane change should be performed if and only if:

314 1. It does not impose an unsafe braking on the target lane following vehicle:

$$\dot{v}_{\text{rear}} \geq -b_{\text{safe}} \quad (36)$$

315 2. It enables the vehicle and (with a politeness factor p) its following vehicles on both current and target
316 lanes to increase their overall acceleration:

$$\Delta \dot{v} + p(\Delta \dot{v}_{\text{rear, current}} + \Delta \dot{v}_{\text{rear, target}}) \geq a_{\text{min}} \quad (37)$$

317 This model describes changes in the target lane L .

318 C Interval Predictor

319 In this section, we design an interval predictor for our system.

320 C.1 Notations

321 For any real variable z , we denote an interval containing z as $\square z = [\underline{z}, \bar{z}]$, such that $\underline{z} \leq z \leq \bar{z}$. As elements of \mathbb{R}^2 , they can be scaled and offset by scalars. This definition is extended element-wise to vector variables.

322 Then, we define several operators over intervals $\square a = [\underline{a}, \bar{a}]$ and $\square b = [\underline{b}, \bar{b}]$

- 324 • The product operator \times

$$\square a \times \square b = [p(\underline{a})p(\underline{b}) - p(\bar{a})n(\bar{b}) - n(\underline{a})p(\bar{b}) + n(\bar{a})n(\underline{b}), \quad (38)$$

$$p(\bar{a})p(\bar{b}) - p(\underline{a})n(\underline{b}) - n(\bar{a})p(\underline{b}) + n(\underline{a})n(\underline{b})] \quad (39)$$

325 where $p(\cdot)$ and $n(\cdot)$ are the projections onto \mathbb{R}^+ and \mathbb{R}^- , respectively.

- 326 • The difference operator $-$

$$\square a - \square b = [\underline{a} - \bar{b}, \bar{a} - \underline{b}] \quad (40)$$

- 327 • The cosine and sine operators

$$\cos(\square z) = [-1 \text{ if } \underline{z} \leq \pi \leq \bar{z} \text{ else } \min(\cos(\underline{z}), \cos(\bar{z})), \quad (41)$$

$$1 \text{ if } \underline{z} \leq 0 \leq \bar{z} \text{ else } \max(\cos(\underline{z}), \cos(\bar{z}))] \quad (42)$$

$$\sin(\square z) = [-1 \text{ if } \underline{z} \leq -\frac{\pi}{2} \leq \bar{z} \text{ else } \min(\sin(\underline{z}), \sin(\bar{z})), \quad (43)$$

$$1 \text{ if } \underline{z} \leq +\frac{\pi}{2} \leq \bar{z} \text{ else } \max(\sin(\underline{z}), \sin(\bar{z}))] \quad (44)$$

- 328 • The inverse operator $/$ over a positive interval $\square z > 0$

$$1/\square z = [1/\bar{z}, 1/\underline{z}] \quad (45)$$

- 329 • Any other function f is assumed increasing on the interval $\square z$ and is applied coefficient-wise

$$f(\square z) = [f(\underline{z}), f(\bar{z})] \quad (46)$$

330 We start with an initial estimate of the intervals over state variables x_I, y_I, v_I and ψ_I . Typically, we use
 331 zero-width intervals centred on the current state observation. Likewise, any variable z used in place of an interval
 332 corresponds to the zero-width interval $[z, z]$.

333 C.2 Intervals for features

334 We use (25) and (35) respectively to derive intervals for the features ϕ_a and ϕ_b from the intervals over the states.

335 We index the front vehicle intervals with the subscript f

$$\square \phi_a = \begin{bmatrix} v_0 - \square v \\ n(\square v_f - \square v) \\ n(\square x_f - \square x - (d_0 + T\square v)) \end{bmatrix} \quad (47)$$

336 and

$$\square \phi_b = \begin{bmatrix} (1 / \square v) \times (y_L - \square y) \\ \psi_L - \square \psi \end{bmatrix} \quad (48)$$

337 C.3 Intervals for controls

338 The controls intervals are derived from (24) and (33)

$$\square a = \square \theta_a^T \times \square \phi_a \quad (49)$$

$$\square \left(\frac{l}{v} \tan \beta \right) = \square \theta_b^T \times \square \phi_b \quad (50)$$

339 **C.4 Intervals for velocity and heading**

340 The velocity interval is derived from (22) and the heading interval from (23)

$$\square \dot{v} = \square a \quad (51)$$

$$\square \dot{\psi} = \square \left(\frac{l}{v} \tan \beta \right) \quad (52)$$

341 **C.5 Intervals for positions**

342 Likewise, the positions interval are derived from the kinematics (20) and (21)

$$\square \dot{x} = \square v \times \cos(\square \psi) \quad (53)$$

$$\square \dot{y} = \square v \times \sin(\square \psi) \quad (54)$$