Sampling-based approximate optimal control under temporal logic constraints

Jie Fu¹ Ivan Papusha² Ufuk Topcu²

¹Worcester Polytechnic Institute

²University of Texas at Austin

HSCC, 2017

April 18-20, 2017 Fu et al.





- 2 Hybrid system formulation
- Approximate-optimal temporal logic planning



Integrate importance sampling

Challenges: Large-scale systems





Challenges: Stringent and complex constraints WPI



A widely explored approach





- Abstraction results in high-dimensional systems Scalability.
- Seek an alternative: No explicit abstraction!
- no free lunch: must resort to approximations

April 18-20, 2017 Fu et al.



- A dynamical system: $\dot{x} = f(x, u)$, x: state, u: input.
- Performance measured by cost function

$$J(x,u) = \int_0^T \ell(x,u) dt + G(x(T),u(T))$$

- Labeling: Atomic propositions \mathcal{AP}
- Specifications in temporal logic: φ over \mathcal{AP} .
- Goal: Design a controller u : [0, T] → U such that the system minimizes the cost while strictly satisfying φ.

From LTL to Automaton



For a given formula φ (subset of LTL), the automaton is

$$A_{\varphi} = \langle Q, 2^{AP}, \delta, q_0, F \rangle.$$

•
$$Q = \{q_0, q_1, q_2, q_3, q_4\}$$

- $\delta: Q \times 2^{AP} \rightarrow Q$ (labeled directed edges).
- q₀: initial state
- F: accepting states.



Figure: $\varphi = \Diamond (R_1 \land \Diamond (R_2 \land \Diamond R_3)) \land \Box \neg Obs.$



An **LTL co-safety formula** can be represented by a deterministic finite state automaton,

 $\underbrace{\rho=q_0q_1\ldots q_n}_{\text{a finite good prefix}}$ is accepting if $q_n\in F,$

where *F* is the set of accepting states in A_{φ} .

- can contain the X (next), ◊ (eventually), U (strong until).
- negation only occurs in the front of atomic propositions.

Limited expressiveness:

- Reachability: \u00f3goal.
- Safety: ¬obs U goal.
- Sequencing: $(\text{goal}_1 \land (\text{goal}_2 \land (\text{goal}_3)))$.
- cannot specify fairness and recurrence properties.

Given a dynamical system

$$\dot{x} = f(x, u)$$

Find *u*^{*} that minimizes

$$J(x_0, u) = \int_{t=0}^{T} \ell(x, u) dt$$

subject to
$$\exists t_0, t_1, \dots, t_k$$

 $\delta(q_i, x(t_i)) = q(t_{i+1}),$
 $i = 0, 1, \dots, k$
 $q_f \in F, x(t_k) = x_f.$

How do we develop a scalable control method that avoids discretizing the state space?

Approximate optimal control for co-safe LTL

Given a dynamical system

 $\dot{x} = f(x, u)$

Find u* that minimizes

$$J(x_0, u) = \int_{t=0}^{T} \ell(x, u) dt$$

subject to
$$\exists t_0, t_1, \dots, t_k$$

 $\delta(q_i, x(t_i)) = q(t_{i+1}),$
 $i = 0, 1, \dots, k$
 $q_f \in F, x(t_k) = x_f.$





Idea: The **unknown** value function can be **arbitrarily closely** approximated by a linear combination of bases.

$$V_q(x) \simeq \hat{V}_q(x) = \sum_{i=1}^n w_{i,q} \phi_{i,q}(x)$$

e.g. polynomial
$$\phi_q(x)$$
 : 1, x^2, x^3, \ldots

unknown value function—**unknown** (finite many) parameters *w_i*.



VPI

April 18-20, 2017 Fu et al.



$$\begin{split} \min_{u} \left(\frac{\partial V_{q}^{\star}}{\partial x} f(x, u) + \ell(x, u) \right) &= 0, \; \forall q \in Q, \; \forall x \in \mathsf{Inv}(q) \\ V_{q}^{\star}(x) - V_{q'}^{\star}(x) &= s(q, q'), \; q' = \delta(q, L(x)). \end{split}$$

and boundary condition $V_q(x_f) = 0$, $\forall q \in F$.



$$egin{aligned} &rac{\partial \hat{m{V}}_{m{q}}}{\partial x}f(x,u)+\ell(x,u)\geq 0, \quad orall m{q}\in m{Q}, \quad orall x\in {
m Inv}(m{q}), \quad orall u\in \mathcal{U}, \ &\hat{m{V}}_{m{q}}(x)-\hat{m{V}}_{m{q}'}(x)=m{s}(m{q},m{q}'), \quad m{q}'=\delta(m{q},L(x)). \end{aligned}$$

and boundary condition

$$\hat{V}_q(x_f) = 0, \forall q \in F.$$

Approximation \hat{V} is a **lower bound** [CDC'16] on the true value function, i.e., $\hat{V}_{q_0}(x_0) \leq V_{q_0}^{\star}(x_0)$



$$egin{aligned} &rac{\partial \hat{m{V}}_{m{q}}}{\partial x}f(x,u)+\ell(x,u)\geq 0, \quad orall m{q}\in m{Q}, \quad orall x\in {
m Inv}(m{q}), \quad orall u\in \mathcal{U}, \ &\hat{m{V}}_{m{q}}(x)-\hat{m{V}}_{m{q}'}(x)=m{s}(m{q},m{q}'), \quad m{q}'=\delta(m{q},L(x)). \end{aligned}$$

and boundary condition

$$\hat{V}_q(x_f) = 0, \forall q \in F.$$

Approximation \hat{V} is a **lower bound** [CDC'16] on the true value function, i.e., $\hat{V}_{q_0}(x_0) \leq V_{q_0}^{\star}(x_0)$



$$\begin{split} \max_{w} \hat{V}_{q_0}(x_0) \\ \text{subject to:} & \frac{\partial \hat{V}_q}{\partial x} f(x, u) + \ell(x, u) \geq 0, \quad \forall q \in Q, \quad \forall x \in \mathsf{Inv}(q), \\ & q_0 = \delta(q_{\textit{init}}, L(x_0)), \\ & \hat{V}_q(x) - \hat{V}_{q'}(x) = s(q, q'), \quad q' = \delta(q, L(x)). \\ & \hat{V}_q(x_f) = 0, \forall q \in F. \end{split}$$



Given *w* (weight parameter), and ϕ_q (vector of basis functions),

$$\begin{split} \max_{w} w_{q_{0}}^{T} \phi_{q_{0}}(x_{0}) \\ \text{subject to:} & \frac{\partial w_{q}^{T} \phi_{q}}{\partial x} f(x, u) + \ell(x, u) \geq 0, \quad \forall q \in Q, \quad \forall x \in \mathsf{Inv}(q), \\ & q_{0} = \delta(q_{\textit{init}}, L(x_{0})), \\ & w_{q}^{T} \phi_{q}(x) - w_{q'}^{T} \phi_{q'}(x) = s(q, q'), \quad q' = \delta(q, L(x)). \\ & w_{q}^{T} \phi_{q}(x_{f}) = 0, \forall q \in F. \end{split}$$

where $\hat{V}_q(x) = \sum_i w_{i,q} \phi_{i,q}(x)$ as the value function approximation.



 LTI system with quadratically representable invariance regions and guards ⇒ SDP [CDC'16] and SYDAR tool

pip install sydar

- Nonlinear system: Semi-infinite program.
- Existing algorithms [R Hettich, 1993] are inefficient with **non-differentiable** objective functions.

Propose: An importance-sampling based search!



Our method is based on MRAS [Hu, Fu, & Marcus, 2007], a general sampling-based method for global optimization

$$x^{\star} \in \operatorname{argmax}_{x \in \mathcal{X}} H(x), \quad \mathcal{X} \subseteq \mathbb{R}^{n}.$$



Our method is based on MRAS [Hu, Fu, & Marcus, 2007], a general sampling-based method for global optimization

$$x^{\star} \in \operatorname{argmax}_{x \in \mathcal{X}} H(x), \quad \mathcal{X} \subseteq \mathbb{R}^{n}.$$

Define a sequence of reference distributions {g_k(·)} that converges to a "Dirac" around x*, e.g.,

$$g_k(x) = rac{H(x)g_{k-1}(x)}{\int_Z H(x')g_{k-1}(x')
u(dx')}, \quad k = 1, 2, \dots$$



Our method is based on MRAS [Hu, Fu, & Marcus, 2007], a general sampling-based method for global optimization

$$x^{\star} \in \operatorname{argmax}_{x \in \mathcal{X}} H(x), \quad \mathcal{X} \subseteq \mathbb{R}^{n}.$$

Define a sequence of reference distributions {g_k(·)} that converges to a "Dirac" around x*, e.g.,

$$g_k(x) = rac{H(x)g_{k-1}(x)}{\int_Z H(x')g_{k-1}(x')
u(dx')}, \quad k = 1, 2, \dots$$

2 Approximate the exact reference distributions {g_k(·)} by a parameterized family of distributions {p(·, θ) | θ ∈ Θ}.



Our method is based on MRAS [Hu, Fu, & Marcus, 2007], a general sampling-based method for global optimization

$$x^{\star} \in \operatorname{argmax}_{x \in \mathcal{X}} H(x), \quad \mathcal{X} \subseteq \mathbb{R}^{n}.$$

Define a sequence of reference distributions {g_k(·)} that converges to a "Dirac" around x*, e.g.,

$$g_k(x) = rac{H(x)g_{k-1}(x)}{\int_Z H(x')g_{k-1}(x')
u(dx')}, \quad k = 1, 2, \dots$$

- 2 Approximate the exact reference distributions {g_k(·)} by a parameterized family of distributions {p(·, θ) | θ ∈ Θ}.
- Generate a sequence of parameters {θ_k} by minimizing the KL divergence between g_k(·) and p(·, θ),

$$D_{KL}(g_k, p(\cdot, \theta)) := \int_Z \ln \frac{g_k(x)}{p(x, \theta)} g_k(x) \nu(dx).$$

A simple/naive formulation



Add a **terminal cost** to penalize value/policy that does not satisfy the temporal logic constraints.

minimize_W
$$J(x_0, u) + \overline{J} \times (1 - 1_F(q_f)).$$

• Given \hat{V} , derive *u* from the HJB—assuming that \hat{V} is optimal.

$$\begin{split} u(x,q) &= \arg\min_{u \in \mathcal{U}} \{ \frac{\partial \hat{V}_q(x)}{\partial x} \cdot f(x,u) + \ell(x,u) \} \\ \hat{V}(x,q) &= \min_{q'} \{ \hat{V}(x,q') + s(x,q,q') \}, \\ \forall x \in G_e, \ \forall e = (q,\sigma,q') \in E, q(t_k^+) = \delta(q(t_k^-), L(x(t_k^+))). \end{split}$$

 Approximate actor-critic NN fails with local optimality and discontinuity in the value function!

Sampling-based approximate optimal planning WPI

Problem: randomly sampling weights that satisfy all the constraints that are also optimal is a **rare** event.



Figure: The overview of the proposed sampling-based optimal control for LTL constraints.

Learning from "bad" weights

- A "bad" weight results in a policy that violates the constraints.
- Difficulty with a good start.
- **Solution:** Learning from the "bad" weights using guided search. $(R_1 \land \Diamond (R_2 \land \Diamond R_3)) \land \Box \neg Obs.$



- Red: satisfies partial spec.
- Black: satisfies the spec but not optimal.

Both provide information about optimal policy.



Idea:

- Instead of throwing away unsatisfying weights,
- ... penalizes with a cost.

Rank-guided weighting of samples.

- rank(q) = 0 for all final q.
- $ank(q) = \min_{q' \in Q} \{ rank(q') + 1 \mid \exists A \in 2^{AP}, \delta(q, A) = q' \}.$

A state-dependent terminal cost:

$$h(q) = \left\{egin{array}{cc} 0 & ext{for all } q \in F. \end{array}
ight.^{ ext{start}} \cdot \ ext{rank}(q) imes c & ext{otherwise.} \end{array}
ight.$$





Idea:

- Instead of throwing away unsatisfying weights,
- ... penalizes with a cost.

Rank-guided weighting of samples.

- rank(q) = 0 for all final q.
- $one and a (q) = \min_{q' \in Q} \{ \operatorname{rank}(q') + 1 \mid \exists A \in 2^{AP}, \delta(q, A) = q' \}.$

A state-dependent terminal cost:

$$h(q) = \left\{egin{array}{cc} 0 & ext{for all } q \in F. \ ext{rank}(q) imes c & ext{otherwise.} \end{array}
ight.$$





Idea:

- Instead of throwing away unsatisfying weights,
- ... penalizes with a cost.

Rank-guided weighting of samples.

- rank(q) = 0 for all final q.
- $ank(q) = \min_{q' \in Q} \{ rank(q') + 1 \mid \exists A \in 2^{AP}, \delta(q, A) = q' \}.$

A state-dependent terminal cost:

$$h(q) = \left\{egin{array}{cc} 0 & ext{for all } q \in F.\ ext{rank}(q) imes c & ext{otherwise.} \end{array}
ight.$$



For linear quadratic system, the optimal value function is polynomial. Red polynomial basis $\phi(x) = [x_1x_2, x_1^2, x_2^2]$, for each $q \in Q$. The number of total weight vectors 15.





Figure: Automaton \mathcal{A}_{φ_1} for $\varphi_1 = (A \rightarrow \Diamond B) \land (C \rightarrow \Diamond B) \land (\Diamond A \lor \Diamond C)$.

The convergence of weight vectors (converge after 10-12 iterations. Each iteration uses 100 samples.)





JPI

Example: Dubins car optimal control with LTL @WPI

Specification: Visit the target while avoiding the obstacles. A mixture of bases:

• Localized Radial Gaussian basis:

$$\phi(\mathbf{x}) = \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2}\right)$$

for pre-selected discrete centers (in x-y coordinate) x_i and σ .

Trigonometric functions:
 cos(θ), sin(θ).



Example: Dubins car optimal control: LTL co-s

Specification: visit regions A, B, C in any order and avoid the obstacles.

A mixture of bases:

• Localized Radial Gaussian basis:

$$\phi(\mathbf{x}) = \exp\left(\frac{-\|\mathbf{x} - \mathbf{x}_i\|^2}{2\sigma^2}\right)$$

for pre-selected discrete centers (in x-y coordinate) x_i and σ .

Trigonometric functions:
 cos(θ), sin(θ).





- We proposed an importance sampling based approximate optimal control algorithm under temporal logic constraints.
- Introducing rank-guided policy search similar to reward shaping — to enable learning from bad samples.

Current and future work:

- Feature selection: How to select a sparse set of basis function.
- Good starting point matters: Value function relates to control Lyapunov function for switched linear systems.