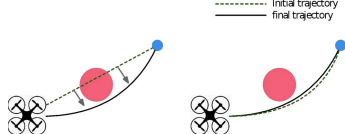




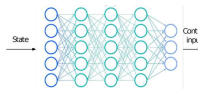
Guided Policy Search using Sequential Convex Programming For Initialization of Trajectory Optimization Algorithms

STUDENTS: Taewan Kim with Purnanand Elango and Behcet Acikmese

Motivation



- Initial trajectory guess is a key input of trajectory optimization algorithms
- It can impact the speed of convergence and reliability of final trajectory
- We generate the initial trajectory guess by training neural net policy



- Once train the policy, we can generate trajectories with different initial points

Problem Formulation

- Optimal control problem. Input "u" comes from neural net.

$$\min_{\theta, x^i(t_k), u^i(t_k)} \sum_{i=1}^N \int_0^{t_f} J(t, x^i(t), u^i(t)) dt$$

$$\text{s.t. } \dot{x}^i(t) = f(t, x^i(t), u^i(t)),$$

$$x^i(t) \in \mathcal{X}(t), \quad u^i(t) \in \mathcal{U}(t),$$

$$s(t, x^i(t), u^i(t)) \leq 0,$$

$$u^i(t) = \pi_\theta(x^i(t_k)), \quad \forall t \in [0, t_f],$$

$$k = 0, \dots, K-1,$$

$$x^i(0) = x_{\text{init}}^i, \quad x^i(t_f) = x_{\text{final}}^i,$$

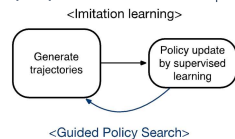
- After optimizing, we can generate a new trajectory by propagating dynamics with the trained policy

$$u(t) = \pi_\theta(x(t_k)) \quad \forall t \in [t_k, t_{k+1}),$$

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

How to train neural net policy?: Guided policy search

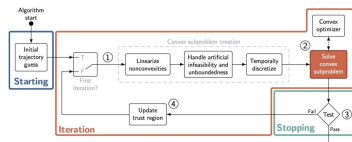
- Imitation Learning (IL)** by trajectory optimization (TO) algorithms
- Generate multiple trajectories by TO as data, then train neural net policy
- However, sometimes the trajectory data is too hard to learn directly
- Guided policy search (GPS)** is motivated from the problem of IL



Levine, Sergey, and Pieter Abbeel. "Learning neural network policies with guided policy search under unknown dynamics." *Advances in neural information processing systems* 27 (2014).

Background: Sequential Convex Programming (SCP)

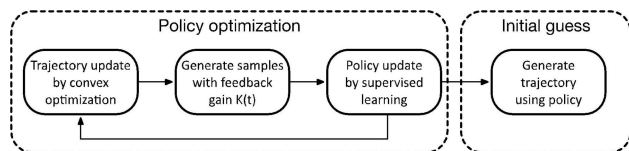
- SCP is trajectory optimization based on convex optimization
- Through linearization and discretization, it approximates the problem to convex sub-problem
- It repeats to this convex optimization problem



[ref] Malaysia, Danylo, et al. "Convex optimization for trajectory generation." *arXiv preprint arXiv:2106.09125* (2021).

Guided Policy Search via SCP

- We develop a new GPS method by employing the idea of SCP
- First, we solve single iteration of SCP
- Second, around the solution of the first step, we generate trajectories with LQR gain
- Third, the trajectory data in second step is used to train the policy
- Repeat the above steps until the convergence of neural net policy



Detail in trajectory update via convex optimization

- Problem formulation

$$\min_{x_k^i, u_k^i, v_k^i} \sum_{i=1}^N \sum_{k=0}^{K-1} J(t_k, x_k^i, u_k^i) + J_{vc}(v) + J_{trp}(u_k^i)$$

$$\text{s.t. } x_{k+1}^i = A_k^i x_k^i + B_k^i u_k^i + z_k^i + v_k^i,$$

$$x_k^i \in \mathcal{X}(t_k), \quad u_k^i \in \mathcal{U}(t_k),$$

$$C_k^i x_k^i + D_k^i u_k^i + r_k^i \leq 0,$$

$$x_0^i = x_{\text{init}}^i, \quad x_{K-1}^i = x_{\text{final}}^i,$$

- Similar with the sub-problem of SCP
- But, it has an additional penalty J_{trp}

$$J_{trp}(u) = w_{trp} \|u_k^i - \pi_\theta(x_k^i)\|_2^2.$$

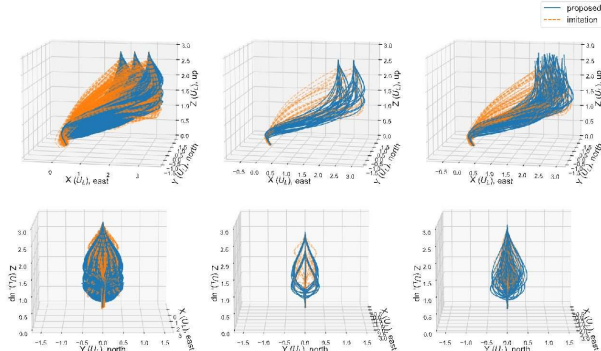
- This penalty enforces the solution is closed to the policy
- So, trajectory update adapts to the neural network

Detail in policy update via supervised learning

- Policy training step
- With data, the policy is trained by supervised learning
- Any neural network optimizer can be used

$$J_p(\theta) = \sum_{s=1}^S \sum_{i=1}^N \sum_{k=0}^{K-1} \|u_{k,s}^i - \pi_\theta(x_{k,s}^i)\|.$$

Application to minimum-fuel powered descent guidance



- 6-dof powered descent guidance problem for a reusable rocket
- Has the following constraints:
- 6-dof rocket dynamics, mass, glide slope angle, angular velocity, tilt angle, thrust, gimbal angle
- Performance comparison

EVALUATION OF COST, CONSTRAINTS AND BOUNDARY CONDITION

Property (Unit)	Validation		Test	
	Proposed	Imitation	Proposed	Imitation
$-\dot{m}(t_f)$ (kg/s)	-1.8591	-1.8648	-1.8610	-1.8670
$\max_k \ c_k\ _\infty$ (s)	0.0031	0.0615	0.0361	0.2112
$\ r_z(t_f) - r_z, f\ _2$ (m)	0.0020	0.0240	0.0081	0.0315
$\ r_x(t_f) - r_x, f\ _2$ (m)	0.0059	0.0507	0.0102	0.0440
$\ h_z(t_f) - h_z, f\ _2$ (°)	0.3074	0.7400	0.5536	0.726
$\ a(t_f) - a_e, f\ _2$ (g)	0.3748	4.9455	0.7976	4.072

INITIALIZATION PERFORMANCE COMPARISON

Property	Proposed	Straight-line
Convergence success rate	100% (100/100)	93% (93/100)
Mean	2.20	8.55
Median	2	5
Standard deviation	0.40	7.90

Future Work, References, and Acknowledgments

- Future works
- Obstacle avoidance problems
- Final time-free problem
- Policy with partial information of state
- Faculty : Behcet Acikmese
- Students : Taewan Kim, Purnanand Elango

