

# Table of Contents

I.Theory.....	2
A.Trajectory.....	2
i.DDP Trajectory Optimization.....	2
ii.Stochastic Trajectory Optimization.....	2
B.MPC with unknown dynamics.....	2
C.Improvement.....	3
i.Line search.....	3
ii.Dual step size of KL divergence constraint.....	3
II. Implement.....	4
A.Design.....	4
B.Point mass.....	4
C.Arm.....	4
III.Result.....	4
A.Box2D.....	4
i.Point mass.....	4
ii.2D link Arm.....	6
IV.Question.....	7
V. Future work.....	8

# MPC\_GPS Demystify and try to apply with unknown dynamics

## I. Theory

### A. Trajectory

- i. DDP Trajectory Optimization
- ii. Stochastic Trajectory Optimization

### B. MPC with unknown dynamics

1. Input: offline trajectory distribution  $p(u_t|x_t)$ , fitted dynamics  $p(x_{t+1}|x_t)$ , current state  $x_t$ .
2. Goal: follow offline trajectory.
3. Output: MPC trajectory distribution  $q(u_{t'}|x_t)$ ,  $\forall t' \in [t+1, t+H]$ .

#### 4. Problem:

4.1. MPC is online trajectory optimization with 2 main attribute:

- a) Use current state  $x_t$  as a initial state  $\rightarrow$  Feedback control.
- b) Short horizon  $\rightarrow$  Need to re-optimize many time  $\rightarrow$  Online optimization.

4.2. Surrogate cost criteria from *MPC\_GPS paper*:

- a) Encourage visit states have high probability from offline LQG solution  
 $p(\tau) = p(u_t|x_t) * p(x_{t+1}|x_t)$ .
- b) Produce good long-horizon behavior while it is short-horizon.
- c) Behavior close to neural network policy  $\pi_\theta(u_t|x_t)$ .

#### 5. Mathematics solution:

5.1. From 4.1 a) I need to have initial state is current state  $x_t$ .

Compute  $p(x_{t'}|x_t)$  by forward pass **with known**  $x_t$  using  $p(u_{t'}|x_{t'})$ ,  $p(x_{t'+1}|x_{t'})$ .

5.2. From 4.2 a & b  $\rightarrow$  The meaningful choice (best from my knowledge) of cost is: cost low (better) when  $p(x_{t'}|x_t)$  high  $\rightarrow -\log p(x_{t'}|x_t)$  is chosen.

5.3. From 4.2 c  $\rightarrow$  MPC solution need to close to neural network policy  $\pi_\theta(u_t|x_t)$ .

But for now, I just want to use MPC as a trajectory optimization, not as a demonstrator from MPC\_GPS. So instead of  $\pi_\theta(u_t|x_t)$ , I use  $p(u_t|x_t)$  - offline LQG solution.

Note: I tried to use  $-\log p(x_{t'}|x_t)$  alone, but it doesn't work.

5.4. The final surrogate cost become:  $\tilde{l}(x_{t'}, u_{t'}) = -\log p(x_{t'}|x_t) - \log p(u_{t'}|x_{t'})$ .

Note: I get rid of  $\nu$ , and don't use  $\eta$  as LQG instead. Check my question 2 below for this.

5.5. The final problem is:  $\min_{q_{ij}(u_t|x_t)} E_{p(x_t, u_t)} [\tilde{l}(x_t, u_t)] - \mathcal{H}(q_{ij}(u_t|x_t))$

5.6. The solution of this problem is the same as LQG,

$$q_{ij}(u_t|x_t) = \mathcal{N}(\tilde{K}_{tij}x_t + \tilde{k}_{tij}, \tilde{Q}_{u,utij}^{-1}).$$

a) To compute  $\tilde{K}$ ,  $\tilde{k}$ , we need to compute gradient and Hessian of  $\tilde{l}(x_t, u_t)$ .

$\frac{\partial}{\partial x}(-\log p(x_{t'}|x_t)) = \Sigma_{t'}^{-1}(x_{t'} - \mu_{t'})$ . In this case  $x_{t'}$  = forward pass **with unknown**  $x_t$ , and use  $x_0$  as initial state using  $p(u_{t'}|x_{t'})$ ,  $p(x_{t'+1}|x_{t'})$ .

b)  $\frac{\partial^2}{\partial x^2}(-\log p(x_{t'}|x_t)) = \Sigma_{t'}^{-1}$ .

$\frac{\partial}{\partial x}(-\log p(u_{t'}|x_t))$ ,  $\frac{\partial^2}{\partial x^2}(-\log p(u_{t'}|x_t))$  are known.

## C. Improvement

### i. Line search

1. Bracket Line search:

1.1. Use bracket line search for  $\eta$  to find the updated distribution  $p(u_t|x_t)$  that KL divergence satisfy the constrained  $KL(p(u_t|x_t)||\hat{p}(u_t|x_t)) \leq \epsilon$ .

1.2. Assume  $p(x_t) = \mathcal{N}(\mu, \Sigma)$ ,  $p(u_t|x_t) = \mathcal{N}(\mu_{u1}, A)$ ,  $\hat{p}(u_t|x_t) = \mathcal{N}(\mu_{u0}, \hat{A})$

1.3.

$$KL = E_{p(x_t)} \left[ \log \frac{p(u_t|x_t)}{\hat{p}(u_t|x_t)} \right]$$

$$= \frac{1}{2} E_{p(x_t)} \left[ -\log |A| + (x - \mu_{u1})^T A^{-1} (x - \mu_{u1}) + \log |\hat{A}| - (x - \mu_{u0})^T \hat{A}^{-1} (x - \mu_{u0}) \right]$$

We have:  $\{(x - \mu_u)^T A^{-1} (x - \mu_u)\} \approx \frac{1}{2} x^T M x + x v + c$

Where  $M = \frac{\partial^2 \{(x - \mu_u)^T A^{-1} (x - \mu_u)\}}{\partial^2 x}$ ,  $v = \frac{\partial \{(x - \mu_u)^T A^{-1} (x - \mu_u)\}}{\partial x}$

$$\text{Then, } KL = \frac{1}{2} E_{p(x_t)} \left[ -\log |A| + \frac{1}{2} x^T M x + x v + c + \log |\hat{A}| - \frac{1}{2} x^T \hat{M} x - x \hat{v} - \hat{c} \right]$$

$$= \left\{ \frac{1}{2} \mu (M - \hat{M}) \mu + \frac{1}{2} \text{Tr}(\Sigma(M - \hat{M})) + \mu(v - \hat{v}) + c - \hat{c} \right\} + \frac{1}{2} \log \frac{|\hat{A}|}{|A|}$$

2. In MPC trajectory optimization, because of unconstrained problem, so I can't use the procedure above to make  $Q_{uu} \succeq 0$  (always positive definite).

2.1. So from *Synthesis and Stabilization of Complex Behaviors through Online Trajectory Optimization 2012*, I use the old regularization trick:  $\tilde{Q}_{uu} = Q_{uu} + \eta \mathcal{I}_m$

2.2.  $\eta$  is updated using bracket line search, but  $\eta_0 = 0$  (original is 1, and use  $\eta$  from previous iteration)

**Question:** Does it work if  $\eta = 10^{-4}$  (small value), because to make  $Q_{uu}$  not singular matrix (full rank ??).

## ii. Dual step\_size of KL divergence constraint

1. From *Learning Contact-Rich Manipulation Skills with Guided Policy Search 2015*, the  $\epsilon$  step size update by using model improvement.
2. Intuition: a larger cost the further we deviate from the previous policy (from *Guided Policy Search as Approximate Mirror Descent 2016*).

The reasonable  $\epsilon$  is the  $\epsilon$  that maximize actual cost improvement.

- 2.1. To do that, we must model the actual cost =  $f(\epsilon) = a\epsilon^2 + b\epsilon$ .

$$actual\_cost = predicted\_cost + noise \Rightarrow b = \frac{predicted\_cost}{\epsilon}.$$

From that can find  $a$  when  $b$ , actual cost are known.

- 2.2. Final problem:  $\max_{\epsilon} f(\epsilon) = a\epsilon^2 + b\epsilon \Rightarrow \epsilon = \frac{-b}{2a}$ .

- 2.3. To pick new  $\epsilon'$ ,  $b = \frac{predicted}{\epsilon}$  where  $\epsilon$  from previous iteration.

- 2.4. Note: From assumption “a larger cost the further we deviate from the previous policy”:

$$KL = \frac{1}{2} \left[ \log \frac{|\Sigma_2|}{|\Sigma_1|} - d + \text{tr}(\Sigma_2^{-1}\Sigma_1) + (\mu_2 - \mu_1)^T \Sigma_2^{-1}(\mu_2 - \mu_1) \right] \text{ (it linear in}$$

covariance  $\Sigma$ , quadratic in mean  $\mu$ ).

$actual\_cost = l_k^k - l_{k-1}^{k-1}$ , and its  $E[actual\_cost]$  is computed below, has the same behavior (linear in covariance  $\Sigma$ , quadratic in mean  $\mu$ ).

1. To compute (Laplace approximation?? - Check Question 3 below) predicted, actual cost of  $E_{p(\tau)}[l(\tau)]$ .

$$1.1. l(\tau) \approx l(x, u) + \begin{bmatrix} x \\ u \end{bmatrix}^T \frac{\partial l(\tau)}{\partial xu} + \frac{1}{2} \begin{bmatrix} x \\ u \end{bmatrix}^T \frac{\partial^2 l(\tau)}{\partial xu, xu} \begin{bmatrix} x \\ u \end{bmatrix}.$$

$$1.2. E[x^T Ax] = \text{Tr}(A\Sigma) + \mu^T A\mu \text{ (From Matrix cookbook).}$$

## II. Implement

### A. Design

### B. Point mass

- i. Add obstacle

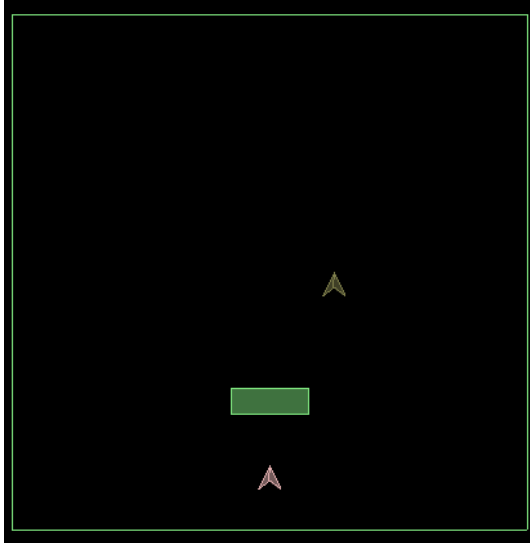
## C. Arm

## III. Result

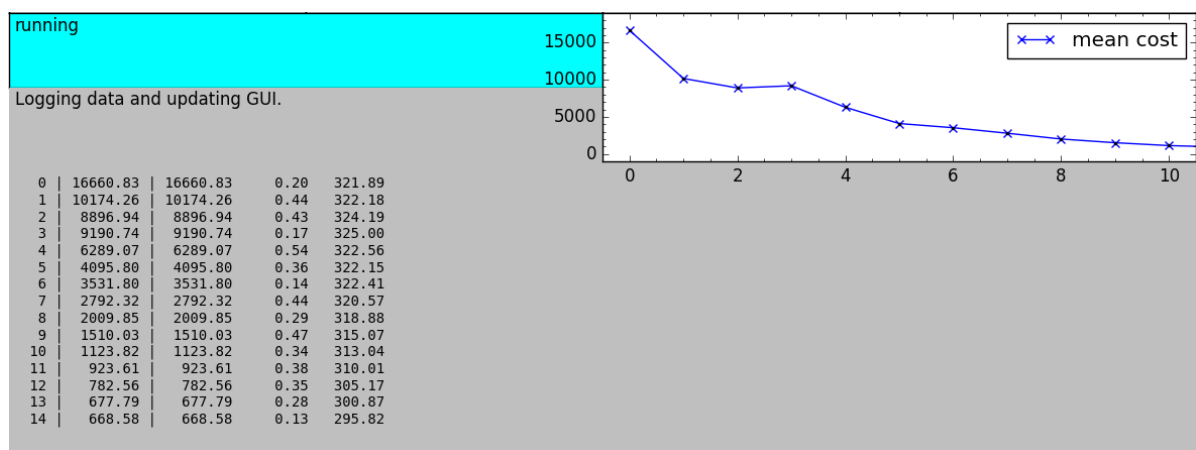
### A. Box2D

#### i. Point mass

1. Offline LQG only (The different from original is that the obstacle)



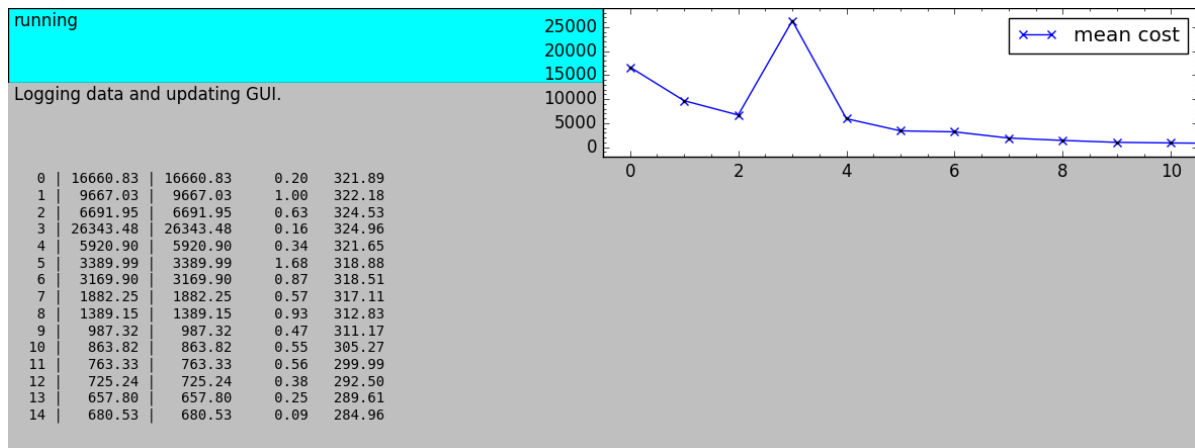
- 1.1. Setup: 15 iterations, 1 obstacle, 1 condition (initial state), 5 samples, T: 100



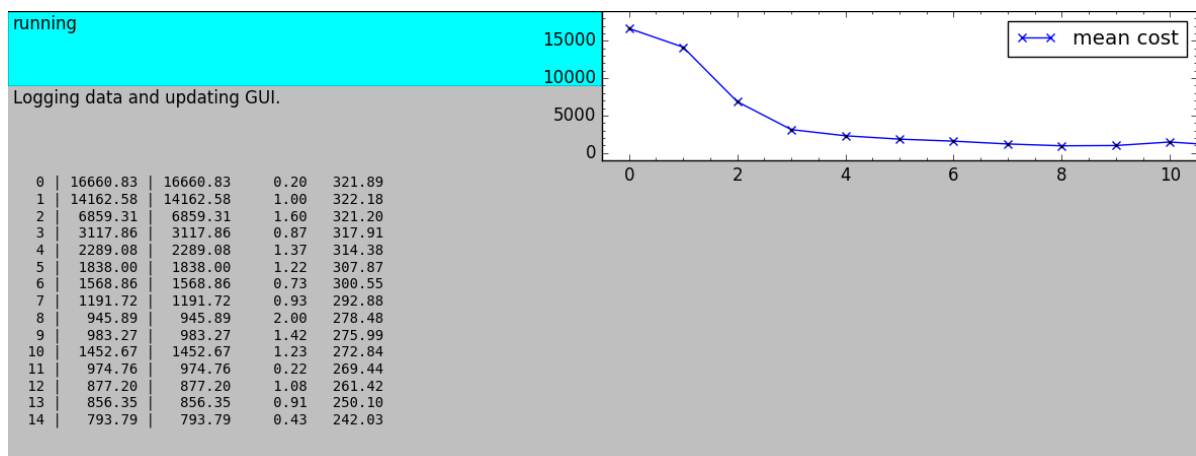
- 1.2. Above image is the result: Row from left to right: iteration, avg\_cost, cost, step, entropy.

2. MPC unconstrained (Check Question 1 below)

- 2.1. Setup: 15 iterations, 1 obstacle, 1 condition (initial state), 5 samples, T: 100



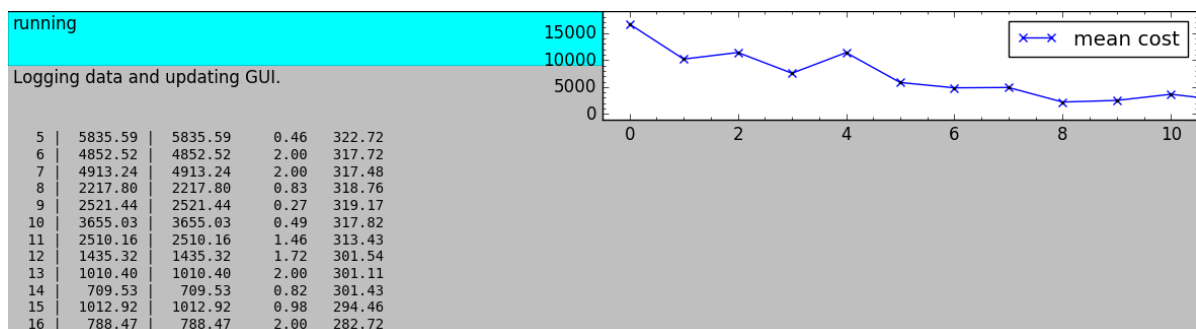
2.2. Short horizon: 10



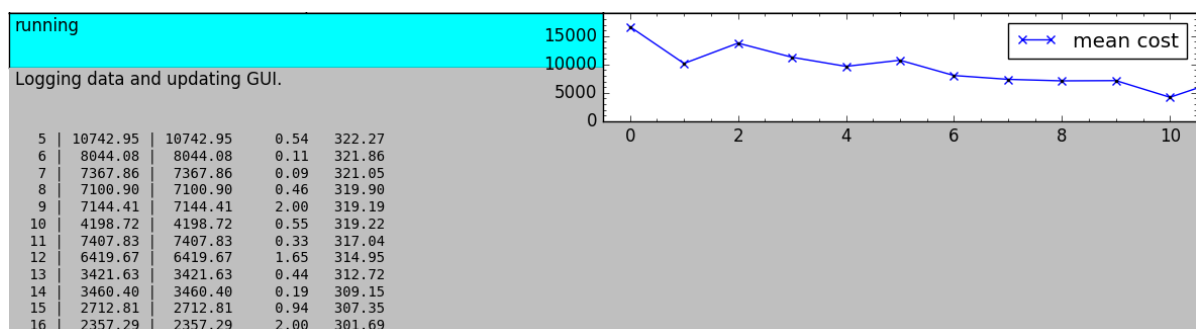
2.3. Short horizon: 5

### 3. MPC constrained (solved by DGD)

3.1. Setup: same as MPC unconstrained



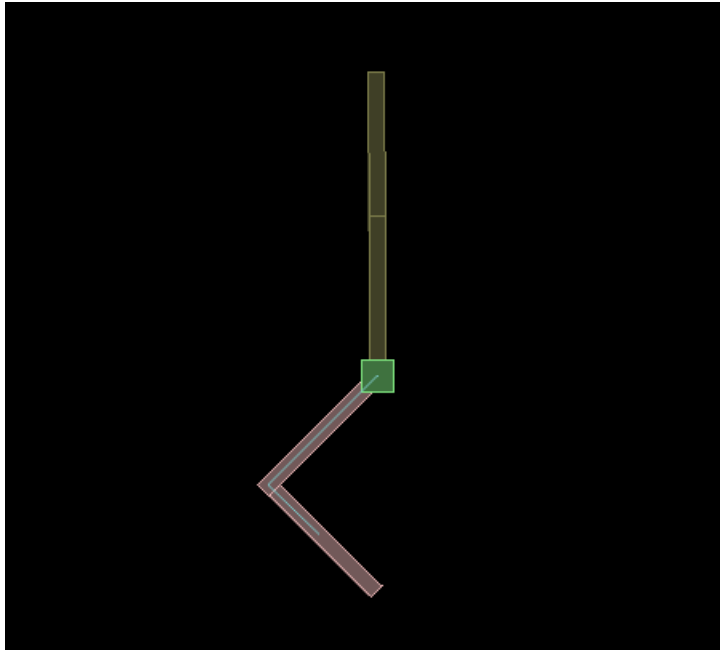
3.2. Short horizon: 10. It work but the cost quite juggle.



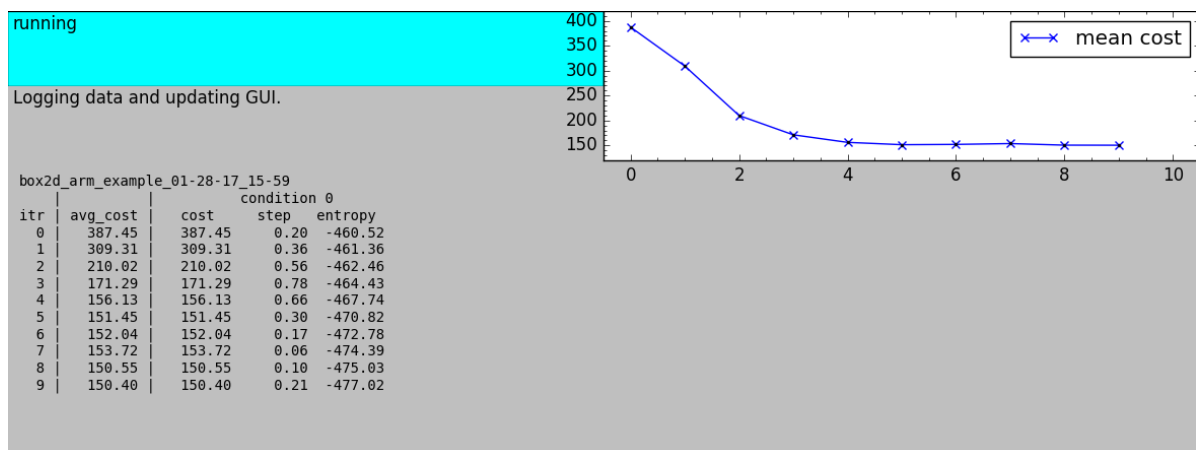
3.3. Short horizon: 5. It cost still high at iteration 16.

## ii. 2D link Arm

1. Offline LQG only (exactly the same as origin Arm world).



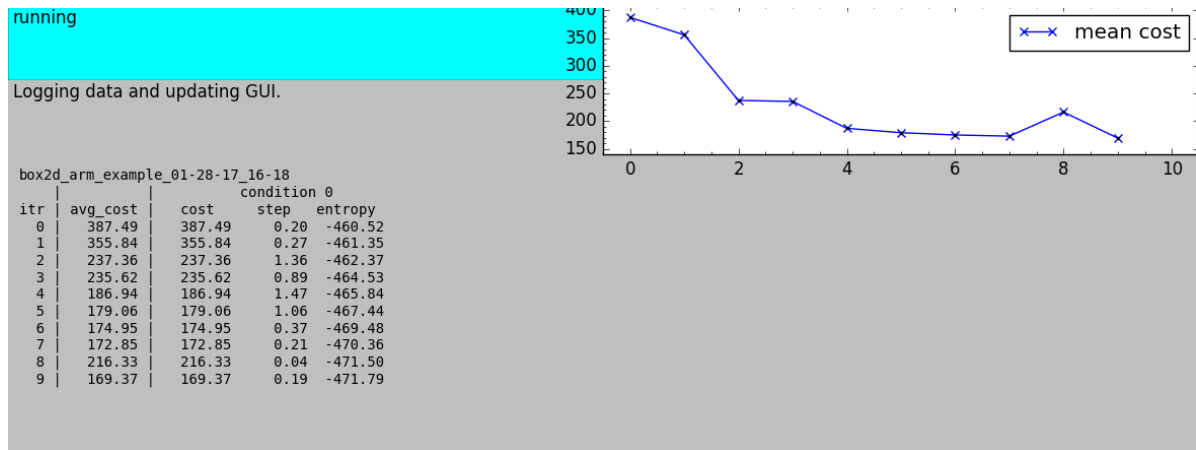
1.1. Setup: 10 iterations, 1 condition (initial state), 5 samples, T: 100



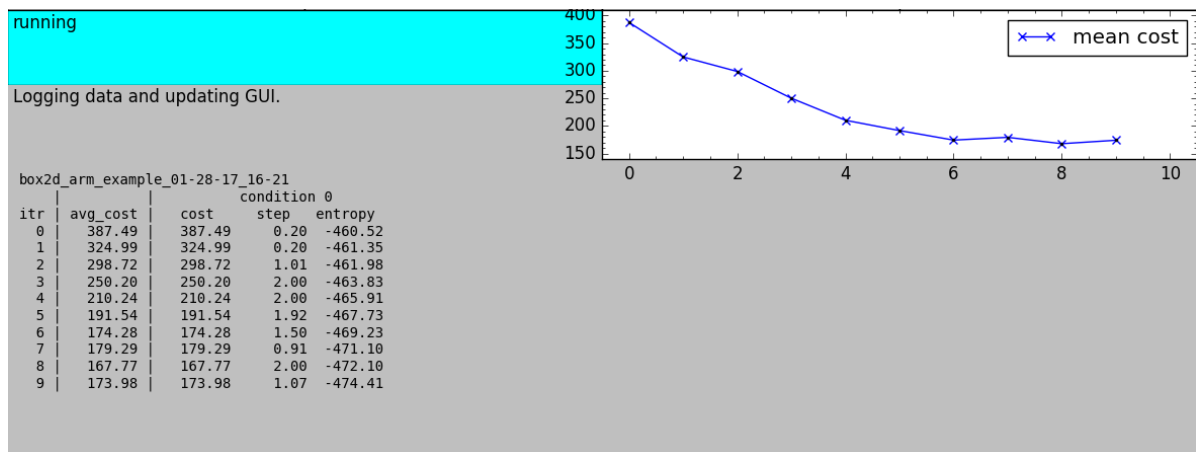
1.2. Above image is the result: Row from left to right: iteration, avg\_cost, cost, step, entropy.

2. MPC unconstrained

2.1. Setup: 10 iterations, 1 condition (initial state), 5 samples, T: 100



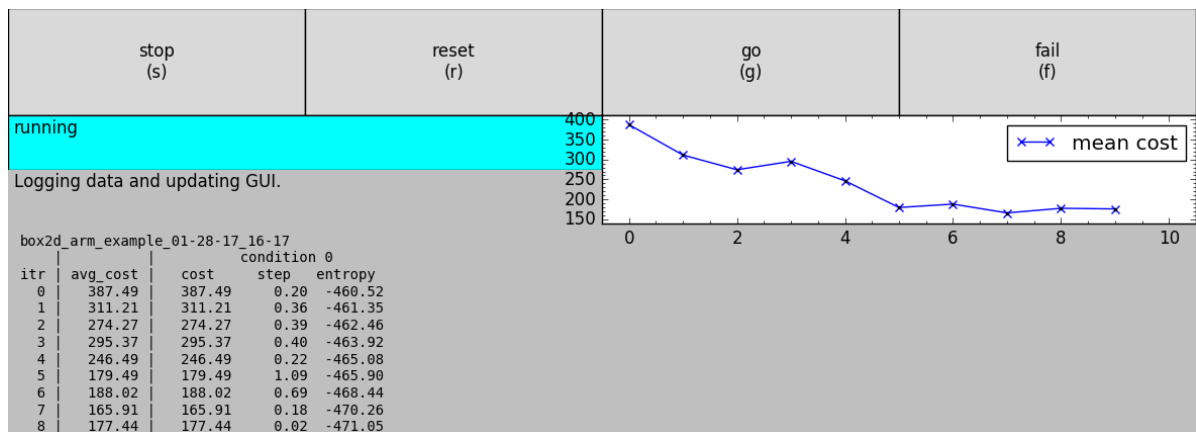
2.2. Short horizon: 10



2.3. Short horizon: 5

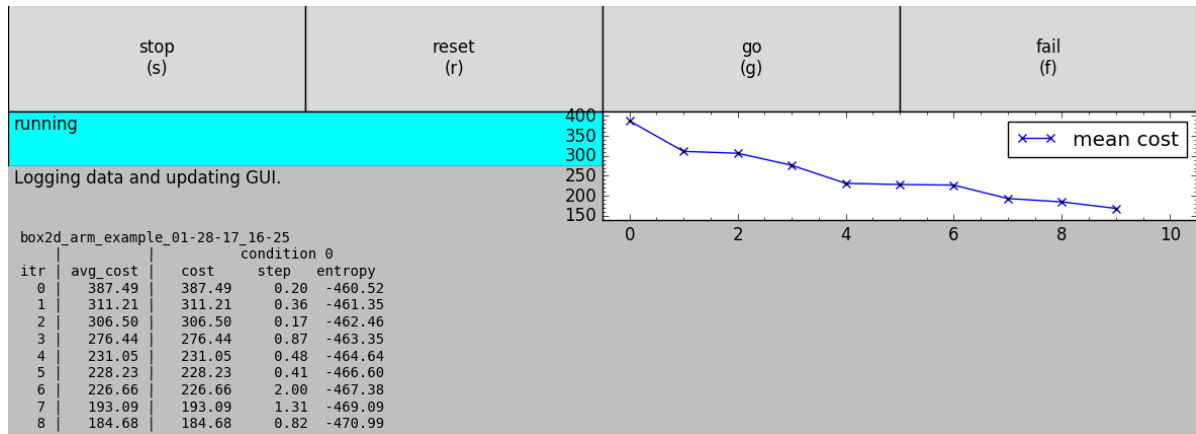
### 3. MPC constrained (solved by DGD)

3.1. Setup: same as MPC unconstrained



3.2. Short horizon: 10





3.3. Short horizon: 5

## IV. Question

1. Why it work when using surrogate cost  $l = -\log q(x_{t'}|x_t) - \log p(u_{t'}|x_{t'})$  without DGD?

**Draft Answer:** Because if I added  $-\log p(u_{t'}|x_{t'})$ . The problem become:

$$\min_{q(u_{t'}|x_{t'})} E_{p(x_t)}[-\log p(x_{t'}|x_t)] + D_{KL}(q(u_{t'}|x_{t'})||p(u_{t'}|x_{t'})) \text{ (KL divergence between}$$

MPC policy and offline LQG policy), it is unconstrained problem, and two term is use log  $\rightarrow$  has the same scale, so  $D_{KL}$  has weight = 1.0 is fine (Just a prediction – not sure).

2. If I mimic exactly the same as offline LQG, then the problem become:

$$\min_{q(u_{t'}|x_{t'})} E_{p(x_t)}[-\log p(x_{t'}|x_t)] + \eta D_{KL}(q(u_{t'}|x_{t'})||p(u_{t'}|x_{t'})) - \eta \epsilon$$

It equivalent to a constraint problem  $D_{KL}(q(u_{t'}|x_{t'})||p(u_{t'}|x_{t'})) \leq \epsilon$ , it work but the cost is not lower (better) than above procedure?

Note: To compute predicted cost and actual cost (Check Dual step\_size of KL divergence constraint) I used the same procedure and cost is raw cost from offline LQG, NOT SURROGATE cost (I tried to use surrogate cost and estimate expected previous cost  $l_{k-1}^{k-1}$  is too far from computed previous cost).

**Draft Answer:** Is this too hard constraint or my way to compute predicted and actual cost to adjust step size is wrong in this case ??

3. Why in the source code, it said Laplace approximation? I though Laplace approximation use to estimate distribution  $p(z)$  with a Gaussian of distribution,  $p(z) \approx q(z) = \mathcal{N}(z_0, A^{-1})$ ,

where  $A = \frac{\partial^2}{\partial z^2} \ln p(z)$  is Hessian around  $z_0$  of  $\ln p(z)$ .

## V. Future work

1. In the box2d\_pointmass world, test more obstacle weight of cost.
2. Demystify BADMM for training neural network policy.
3. Read more on Information Theory.