

---

---

# Control with Reinforcement Learning

---

---

---

# Contents

Comparison of RL and Supervised Learning

Brief introductions of terms and background

Parameter Optimization for Policy (using  $PI^2$ )

Skipping Parameter Optimization with MPC

---

---

# Limits of Supervised Learning (for Robot Control)

Assumes training data *perfectly represents* the variety of data experienced at run-time.

Does not have mechanism to incorporate new data outset of training.

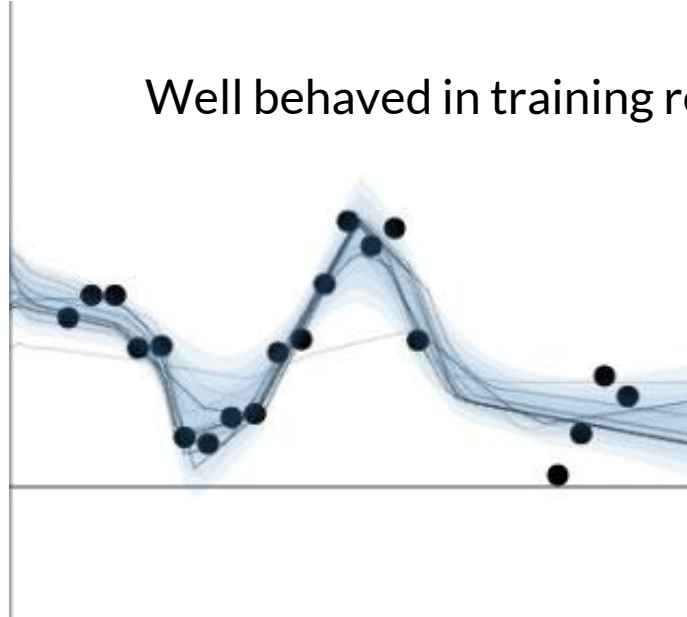
Will probably have unexpected / dangerous behavior when operating out of spec.

---

---

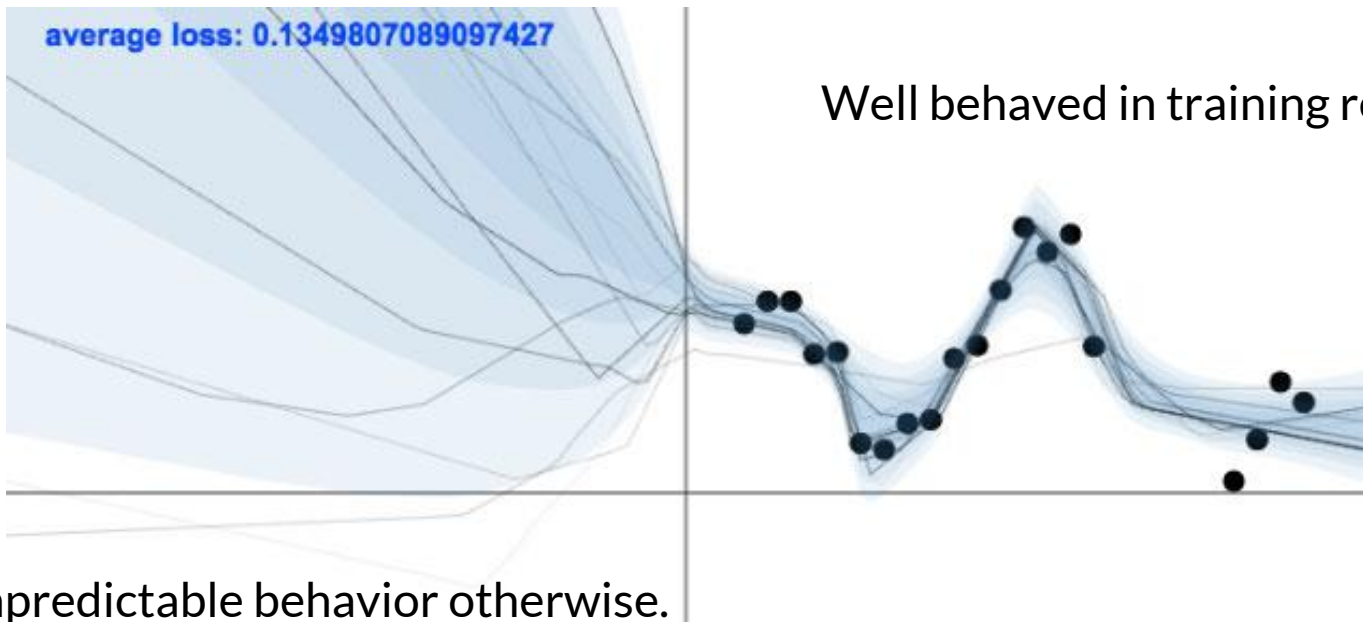
---

Well behaved in training regime...



---

•



---

# How does RL help?

The learner is an *agent*, that takes actions to achieve a high reward / low cost.

Actions can take the agent to new states (new data), that will be incorporated into the agent's behavior.

RL can be viewed as a superset of Supervised Learning.

---

---

# Some Terms

*Trajectory* : A sequence of states and actions, time indexed, where the state at some time  $t$ , depends on the state and action taken at time  $t-1$ .

*Cost function* : a function that takes in a state-action pair, and returns a value denoting 'good', 'bad', or how well that pair accomplishes some goal.

---

# Some Terms in RL

$V(s)$  : Value function. Represents the *value* of being in some state 'x' relative to the current task.

$$V(x_t) = R(x_t) + V(x_{t+1})$$

$Q(s,a)$  : Q-function. Represents the *value* of being in a state 'x' and taking an action 'u'.

$$Q(x_t, u_t) = R(x_t) + Q(x_{t+1}, \pi(x_{t+1}))$$

$a=\pi(s)$  : Policy function. Returns an action for a given state. The optimal policy is the highest scoring action in the Q-function, given an optimal Q function.

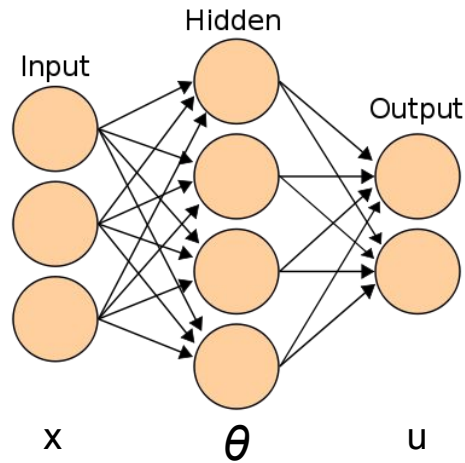
0.22	↖	↖	↖	0.30	↖	0.34	↖	0.38	↖	0.34	↖	0.34	↖	0.38	↖
0.25	→	→	→	0.34	→	0.38	→	0.42	→	0.38	→	0.34	→	0.38	→
0.21	↑					0.46	↑							0.46	↑
0.20	↘	↘	↘	-0.78	↘	0.52	↘	0.57	↘	0.64	↘	0.57	↘	0.52	↘
0.22	↖	↖	↖	0.25	↖	0.27	↖	0.08	↖	-0.36	↖	0.71	↖	0.64	↖
0.25	↖	↖	↖	0.27	↖	0.25	↖	1.20	↖	0.08	↖	0.79	↖	-0.29	↖
0.27	↖	↖	↖	0.30	↖	0.30	↖	1.04	↖	0.97	↖	0.87	↖	-0.21	↖
0.31	↖	↖	↖	0.38	↖	-0.58	↖	-0.44	↖	-0.14	↖	0.74	↖	0.71	↖
0.34	→	→	→	0.42	→	0.46	→	0.52	→	0.57	→	0.64	→	0.71	→
0.31	↘	↘	↘	0.42	↘	0.46	↘	0.52	↘	0.57	↘	0.64	↘	0.71	↘



---

# Function Approx. of Policy

$u = \pi_{\theta}(x)$  : here we say policy  $\pi$  is parameterized by  $\theta$ .



---

# Simple Policy Improvement

*Cross-Entropy Method (one iteration)*

$$\boldsymbol{\theta}_{k=1\dots K} \sim \mathcal{N}(\boldsymbol{\theta}, \Sigma) \quad \text{sample (1)}$$

$$J_k = J(\boldsymbol{\theta}_k) \quad \text{eval. (2)}$$

$$\boldsymbol{\theta}_{k=1\dots K} \leftarrow \text{sort } \boldsymbol{\theta}_{k=1\dots K} \text{ w.r.t } J_{k=1\dots K} \quad \text{sort (3)}$$

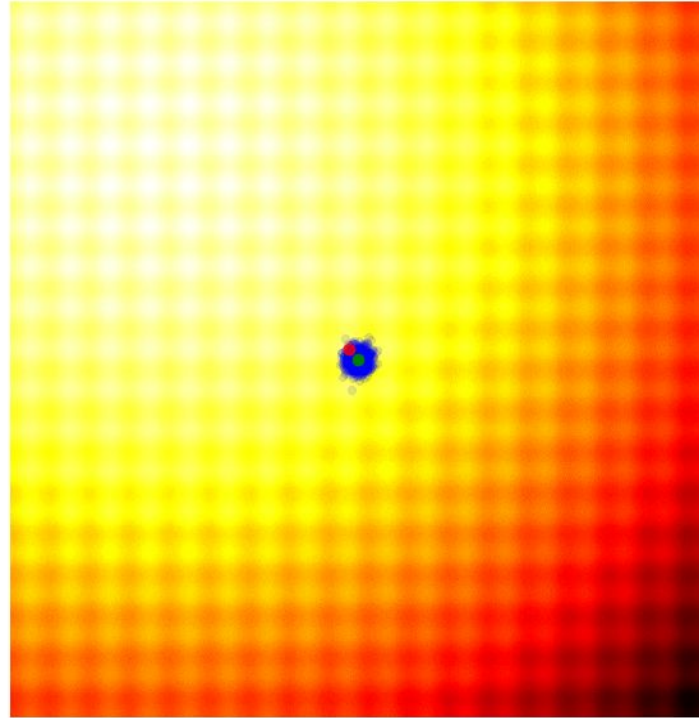
$$\boldsymbol{\theta}^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} \boldsymbol{\theta}_k \quad \text{update (4)}$$

$$\Sigma^{new} = \sum_{k=1}^{K_e} \frac{1}{K_e} (\boldsymbol{\theta}_k - \boldsymbol{\theta})(\boldsymbol{\theta}_k - \boldsymbol{\theta})^\top \quad \text{update (5)}$$

1. Sample  $K$  policies from a distribution.
  2. Evaluate the *1-step score* of each policy with cost function  $J(\boldsymbol{\theta})$
  3. Rank each policy  $\boldsymbol{\theta}_k$  based on its score
  4. Make a new policy; here we average the  $K_e$  best performing policies.
  5. Can also update other parameters with this scoring.
-

---

# CMA-ES



# PI<sup>2</sup>

1. Sample  $K$  policies from a distribution.
2. *Rollout* each policy for  $N$  timesteps, and store states and actions.
3. Evaluate the *performance* of each policy with cost function  $J(\theta)$
4. Weight each trajectory's score against all others, and use them to make a weighted average of our  $K$  policies.

PI<sup>2</sup>

for  $k = 1 \dots K$  **do**

|  $\theta_{k,i=1\dots N} \sim \mathcal{N}(\theta, \Sigma)$

|  $\tau_{k,i=1\dots N} = \text{executepolicy}(\theta_{k,i=1\dots N})$

for  $i = 1 \dots N$  **do**

| for  $k = 1 \dots K$  **do**

| |  $S_{k,i} \equiv S(\tau_{k,i}) = \sum_{j=i}^N J(\tau_{j,k})$

| |  $P_{k,i} = \frac{e^{-\frac{1}{\lambda} S_{k,i}}}{\sum_{k=1}^K [e^{-\frac{1}{\lambda} S_{k,i}}]}$

|  $\theta_i^{new} = \sum_{k=1}^K P_{k,i} \theta_k$

|  $\Sigma_i^{new} = \sum_{k=1}^K P_k (\theta_{k,i} - \theta)(\theta_{k,i} - \theta)^\top$

$\theta^{new} = \frac{\sum_{i=0}^N (N-i) \theta_i^{new}}{\sum_{l=0}^N (N-l)}$

$\Sigma^{new} = \frac{\sum_{i=0}^N (N-i) \Sigma_i^{new}}{\sum_{l=0}^N (N-l)}$

---

# Limitations of $PI^2$

$K$  needs to be somewhat large; takes a lot of compute to do many rollouts.

$\theta$  may also have many parameters, contributing to compute issues.

Trajectory optimization is a hard problem; difficult to find a path from start to finish effectively.

*Do we even need a policy?*

---

---

# MPPI

MPC : Instead of a long, perfect trajectory, we constantly recalculate short trajectories whilst simultaneously using the results. Can make up for modelling errors, run-time perturbations, etc.

Given that a policy & model produces a sequence of controls; can we just optimize a sequence of controls?

```
for t=1:T
    u_t =  $\pi(x_t)$ 
    x_{t+1} = F(x_t, u_t)
end
```

---

---

**Algorithm 2:** MPPI

---

**Given:**  $\mathbf{F}$ : Transition Model;  
 $K$ : Number of samples;  
 $T$ : Number of timesteps;  
 $(\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{T-1})$ : Initial control sequence;  
 $\Sigma, \phi, q, \lambda$ : Control hyper-parameters;  
**while** *task not completed* **do**  
     $\mathbf{x}_0 \leftarrow \text{GetStateEstimate}()$ ;  
    **for**  $k \leftarrow 0$  **to**  $K - 1$  **do**  
         $\mathbf{x} \leftarrow \mathbf{x}_0$ ;  
        Sample  $\mathcal{E}^k = \{\epsilon_0^k, \epsilon_1^k, \dots, \epsilon_{T-1}^k\}$ ;  
        **for**  $t \leftarrow 1$  **to**  $T$  **do**  
             $\mathbf{x}_t \leftarrow \mathbf{F}(\mathbf{x}_{t-1}, \mathbf{u}_{t-1} + \epsilon_{t-1}^k)$ ;  
             $S(\mathcal{E}^k) += \mathbf{q}(\mathbf{x}_t) + \lambda \mathbf{u}_{t-1}^\top \Sigma^{-1} \epsilon_{t-1}^k$ ;  
         $S(\mathcal{E}^k) += \phi(\mathbf{x}_T)$ ;  
     $\beta \leftarrow \min_k [S(\mathcal{E}^k)]$ ;  
     $\eta \leftarrow \sum_{k=0}^{K-1} \exp(-\frac{1}{\lambda}(S(\mathcal{E}^k) - \beta))$ ;  
    **for**  $k \leftarrow 0$  **to**  $K - 1$  **do**  
         $w(\mathcal{E}^k) \leftarrow \frac{1}{\eta} \exp(-\frac{1}{\lambda}(S(\mathcal{E}^k) - \beta))$ ;  
    **for**  $t \leftarrow 0$  **to**  $T - 1$  **do**  
         $\mathbf{u}_t += \sum_{k=1}^K w(\mathcal{E}^k) \epsilon_t^k$ ;  
    SendToActuators( $\mathbf{u}_0$ );  
    **for**  $t \leftarrow 1$  **to**  $T - 1$  **do**  
         $\mathbf{u}_{t-1} \leftarrow \mathbf{u}_t$ ;  
     $\mathbf{u}_{T-1} \leftarrow \text{Intialize}(\mathbf{u}_{T-1})$ ;

---

---

Instead of a policy with parameters  $\theta$ , we instead just keep a sequence of controls  $\mathbf{U} = [\mathbf{u}_0, \mathbf{u}_1, \mathbf{u}_2, \dots]$ ; instead of  $\theta_k$  we have  $\mathbf{U}_k = [\mathbf{u}_0 + \epsilon_0^k, \mathbf{u}_1 + \epsilon_1^k, \mathbf{u}_2 + \epsilon_2^k, \dots]$ .

We rollout each  $\mathbf{U}_k$  and score each trajectory. Score is  $S(\epsilon^k)$ , and cost function is  $q()$ .

The weighting is the same as PI<sup>2</sup>, with a baseline subtraction to scale things.

The weights give us a new  $\mathbf{U}$  instead of new policy.

---





---

# Implementation Details

Task: Go to Pose. When close enough, consider task done.

Avoid for-loops. Use PyTorch's Tensor math. Look up 'mm', 'bmm', and others...

Chance are you will need to smooth your control outputs; you can smooth before you update, or smooth the values sent to the actuators. When adding noise to controls, you should clip to the values your NN model was trained at.

$T=30$ ;  $K=1000$  at 10hz is easily achievable.

---

---

# Test.

Debug each component in parts. Does your NN correctly 'drive' forward? Backwards? Left and right? Sitting still?

Does each part of MPPI do the expected thing? Set T and K to small values, and look at all the numbers. Convince yourself it's doing the correct thing.

---

---

# Useful References

<http://blog.otoro.net/2017/10/29/visual-evolution-strategies/> GIFs

<https://icml.cc/2012/papers/171.pdf> PI<sup>2</sup> with CMA

<https://www.cc.gatech.edu/~bboots3/files/InformationTheoreticMPC.pdf>

<https://arxiv.org/abs/1707.04540> newest MPPI work (has some tips)

---

