

Recitation 8:

Quiz 3 Review

Overview:

Main Topics:

- All papers not marked “optional” since and including 09/27
- PPO, TRPO
- PETS
- MBPO
- DDPG
- AlphaGo, AlphaGoZero
- Evolutionary Strategies
- BC, DAGGER
- GAIL

Quiz Date: Monday 11/08 During Class

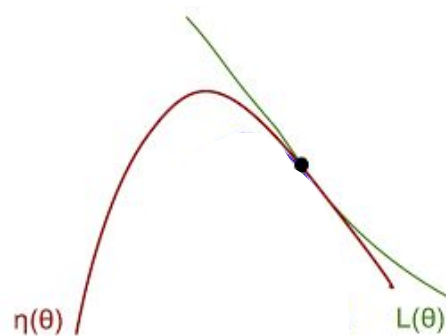
Trust Region Policy Optimization (TRPO) Motivation

$$\text{Goal: } \max_{\theta} \eta(\theta) = \mathbb{E}_{s_0, a_0, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right] = \mathbb{E}_{s \sim \rho_{\pi_{\theta}}} \left[\sum_{t=0}^{\infty} R(s) \right]$$

$$\text{Exact Update: } \eta(\theta_{\text{new}}) = \eta(\theta_{\text{old}}) + \sum_s \rho_{\pi_{\theta_{\text{new}}}}(s) \sum_a \pi_{\theta_{\text{new}}}(a | s) A_{\theta_{\text{old}}}(s, a)$$

Approximation:

$$L_{\theta_{\text{old}}}(\theta_{\text{new}}) = \eta(\theta_{\text{old}}) + \sum_s \rho_{\pi_{\theta_{\text{old}}}}(s) \sum_a \pi_{\theta_{\text{new}}}(a | s) A_{\theta_{\text{old}}}(s, a)$$

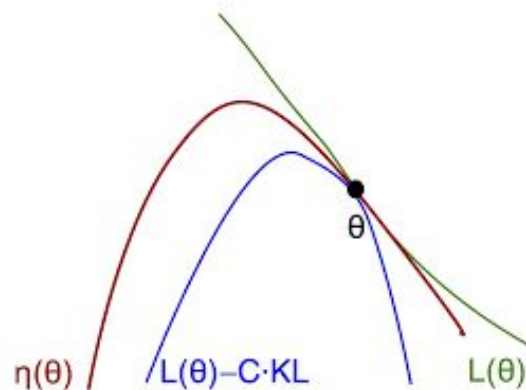


TRPO: monotonic improvement Theorem

We want to construct a lower bound on $\eta(\theta)$

Theorem: $\eta(\theta_{\text{new}}) \geq L_{\theta_{\text{old}}}(\theta_{\text{new}}) - CD_{\text{KL}}^{\text{max}}(\theta_{\text{old}}, \theta_{\text{new}})$

where $C = \frac{4\epsilon\gamma}{(1-\gamma)^2}$, $\epsilon = \max_{s,a} |A_{\theta_{\text{old}}}(s, a)|$



TRPO: In practice

The monotonic improvement Theorem proposes too small steps in practice.

Instead: $\max_{\theta} L_{\theta_{\text{old}}}(\theta)$

subject to $D_{\text{KL}}^{\max}(\theta_{\text{old}}, \theta) \leq \delta$

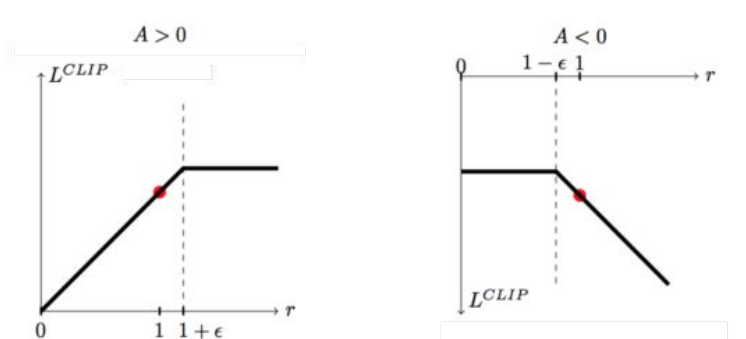
To compute this, TRPO uses the conjugate gradient method, requiring computing the Fisher Information Matrix (FIM), i.e. the the Hessian of the KL divergence.

Proximal Policy Optimization (PPO)

Can we simplify the KL constraint, so we don't have to compute the FIM?

Let the ratio $r_t(\theta) = \pi_\theta(a_t | s_t) / \pi_{\theta_{\text{old}}}(a_t | s_t)$

Define: $\mathcal{L}_{\theta_{\text{old}}}^{\text{CLIP}}(\theta) = \mathbb{E}_{\tau \sim \pi_{\text{old}}} \left[\sum_{t=0}^T \left[\min \left(r_t(\theta) \hat{A}_t^{\pi_{\text{old}}}, \text{clip} \left(r_t(\theta), 1 - \epsilon, 1 + \epsilon \right) \hat{A}_t^{\pi_{\text{old}}} \right) \right] \right]$



Actor-Critic methods (Quick refresher)

Objective: $\max_{\theta} \mathbb{E}_{\tau \sim P_{\theta}} [R(\tau)]$

Gradient update: $\mathbb{E}_{s \sim d^{\pi_{\theta}}(s), a \sim \pi_{\theta}(a|s)} \nabla_{\theta} \log \pi_{\theta}(a | s) [A(s, a; \phi)]$

Deep Deterministic Policy Gradient (DDPG)

What if we instead represent our policy as a deterministic function.

$$\pi : S \mapsto A$$

If our learned Q-function is differentiable, we can directly optimize our policy to maximize the Q-function.

$$\max_{\theta} \mathbb{E}_{\pi_{\theta}} [Q(s, \pi_{\theta}(s))]$$

Gradient update:

$$\nabla_{\theta^{\mu}} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a | \theta^Q) \Big|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^{\mu}} \mu(s | \theta^{\mu}) \Big|_{s_i}$$

To explore, add Gaussian noise during training.

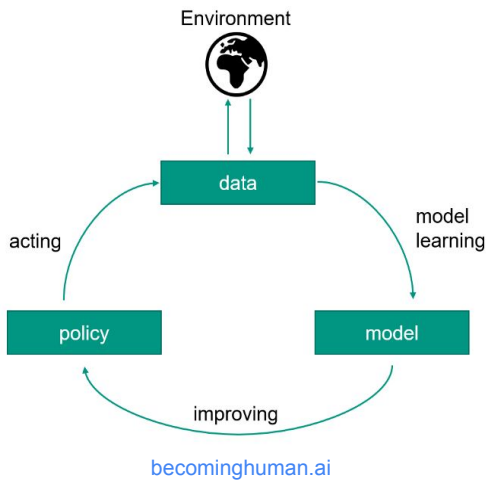
Soft Actor-Critic (SAC)

Increasing the entropy of our policy improves exploration and makes it less likely to get stuck in a local minimum.

SAC adds an entropy reward to encourage higher entropy policies.

$$\max_{\pi} E_{\pi} \left[\sum_t \gamma^t (r(s_t, a_t) + \mathcal{H}_{\pi} [a_t | s_t]) \right]$$

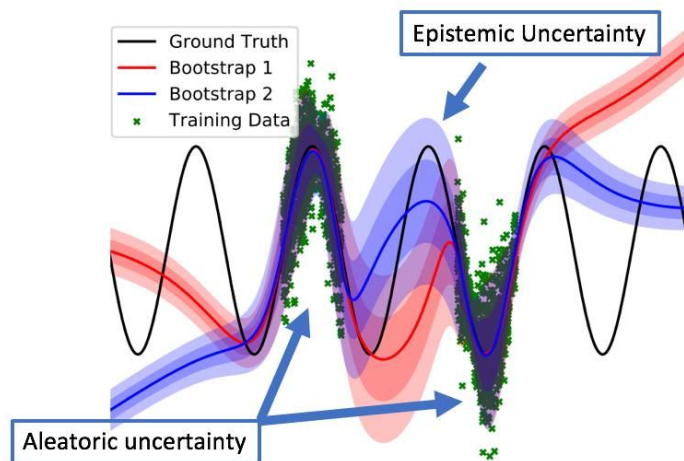
Model-based RL (MBRL)



Main challenge for MBRL: how to deal with model errors?

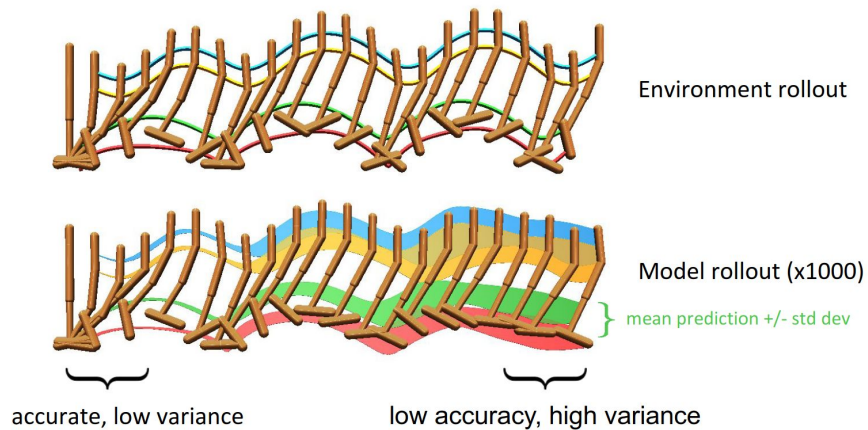
Probabilistic Ensembles with Trajectory Sampling (PETS)

- Uses probabilistic models to estimate Aleatoric uncertainty
- Uses an ensemble of such models to estimate Epistemic uncertainty
- Uses Model Predictive Control (MPC) to compute a new policy at every step:
 - Use particle-based sampling method to estimate returns for a given policy
 - Use CEM to optimize the policy



Model Based Policy Optimization (MBPO)

- Longer trajectories result in larger errors
- By simply decreasing the trajectory length, we can increase accuracy
 - Note: this only works on domains with dense reward functions
- MBPO uses SAC to optimize the policy

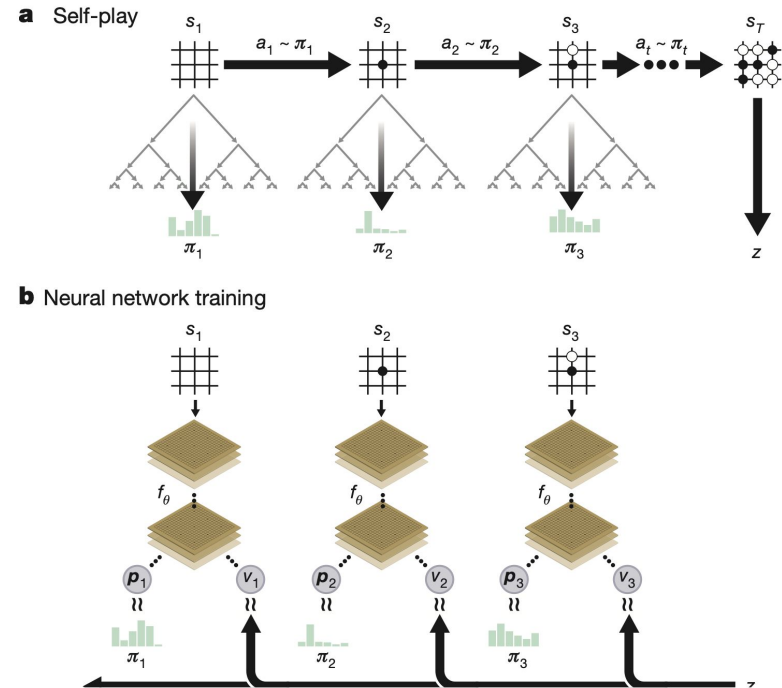


AlphaGo

- Go, massive state-action space
- Use combined policy, value network to prune MCTS tree depth and breadth
- 19x19 board passed through CNN for feature extraction
- Start with SL p_{σ} , p_{π} from expert moves
- Then train p_{ρ} with RL to optimize game outcomes in self-play
 - Need to re-adjust objective to winning game
- Finally train v_{θ} that predicts winners of self-play games

AlphaGo Zero

- Combine policy/value networks into single architecture
- No expert knowledge needed
- Uses MCTS as a policy improvement operator
- Each time step is used to update model parameters using (s_t, π_t, z_t) tuples



Evolutionary Strategies

- Permute → Test mutations → Select the best for the next round → Repeat
- “Black-box” optimization
 - Directly search over policy parameter space
 - No gradient information
 - No reward information
 - No state-structure information
 - May get stuck at local optima
- CEM (Cross Entropy Method)
 - Sample using a multivariate gaussian w/ diagonal covariance (independent dimensions)
 - Update mean and variance using elites
 - Great low-dimensional performance
- CMA-ES (Covariance Matrix Adaptation Evolutionary Strategy)
 - Sample using a multivariate gaussian w/ full covariance (correlated dimensions)
 - Update mean and variance using elites

Behavioral Cloning (BC), DAGGER

- Methods for imitation learning (making a novice act like an expert)
- BC:
 - Use the expert to generate and label random training data
 - Use supervised learning to train the novice
 - Problem: Distribution Shift between Test
- Solution: Data Aggregation (DAGGER)
 - “Fuse” the expert and novice distributions by using the expert to relabel the novice’s own attempts
 - Keep learning from your past knowledge with aggregation
 - No “forgetting”

Generative Adversarial Imitation Learning (GAIL)

- Learn without a reward signal or interaction with an expert
 - Start only with a set of trajectories from the expert, no more querying allowed
- Use a Discriminator network as a cost function to train a parameterized policy
 - $D(s,a) \rightarrow [0,1]$
 - Policy loss: $\log(D_{\text{updated}}(s,a))$, use TRPO for update step
- Outperforms BC

Ultimate objective is to find a saddle point of:

$$\mathbb{E}_{\pi} [\log(D(s, a))] + \mathbb{E}_{\pi_E} [\log(1 - D(s, a))] - \lambda H(\pi)$$

Some Questions...

When does GAIL outperform other IL methods (i.e., BC & DAGGER) and why ?

Some Questions...

When does GAIL outperform other IL methods (i.e., BC & DAGGER) and why ?

GAIL works better when the expert trajectory distribution is multi-modal, because minimizing KL divergence is only suited to approximate uni-modal expert distribution.

Some Questions...

What is maximization bias?

Some Questions...

What is maximization bias?

Overestimation of a value function due to the fact that the same function which is being optimized is used to evaluate its own performance. (Happens on q-value bootstrapping step)

Some Questions...

Describe one of the drawbacks of DDPG.

Some Questions...

Describe one of the drawbacks of DDPG.

DDPG can be unstable and heavily reliant on finding the correct hyper parameters for the current task. This is caused by the algorithm continuously over estimating the Q values of the critic (value) network.