

Lecture 7: Policy Gradient methods

→ optimize the policy directly

[Model free]

Previously, VFA $\rightarrow V_{\theta}(s) \approx V^{\pi}(s)$

$$Q_{\theta}(s, a) \approx Q^{\pi}(s, a).$$

and policy generated based on value function.

Here, parameterize policy $\rightarrow \pi_{\theta}(s, a) = P(a | s, \theta)$.

main motivation

↳ Large scale RL.

↳ some function approximator.

Policy based

(In value based, we learnt a value function)

→ No value function

→ Learnt policy

Use Gradient Ascent to maximize reward.

(Why policy based?)

→ policy can be more compact representation

→ better convergence properties

→ effective in continuous space

(we had 'max' in value function)

→ learn stochastic policies

→ directly follow the policy gradient

↳ can't do in continuous action space.

Disadvantages:

→ Naïve can be inefficient, high variance, slower.

→ Policy can be local optimum.

converged to,

Why stochastic policy?

- deterministic policy \rightarrow can be easily exploitable
- eg. Rock paper scissor \rightarrow if we always do Rock, the opponent will do paper and defeat us every time
- Nash equilibrium \rightarrow game theoretic equivalent of optimality.
- \therefore optimal behaviour \rightarrow stochastic \rightarrow uniform random

\rightarrow state aliasing \rightarrow when markov property doesn't hold
eg. partially observable MDPs \rightarrow see some part - some features that give uncomplete information.

eg. we have features

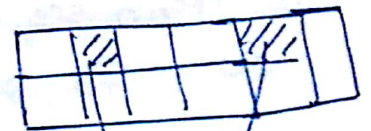
(for eg. in slide)

$$\phi(s, a) = \mathbb{I}(\text{wall to } N, a = \text{move } E).$$

Difference betⁿ value based and policy based RL:

$$Q_{\theta}(s, a) = f(\phi(s, a), \theta)$$

$$\pi_{\theta}(s, a) = g(\phi(s, a), \theta)$$



then deterministic policy will always choose the same method. \therefore can get stuck for a long time

\leftarrow in the case of value based RL (greedy)

But stochastic policy \rightarrow policy to randomly move E/W in marked states

$$\begin{aligned}\pi(\text{move } E) &= 0.5 \\ \pi(\text{move } W) &= 0.5\end{aligned}$$

stochastic $>$ deterministic for state

aliasing \rightarrow not a complete MDP

\rightarrow where deterministic always reaches optimal policy

Policy objective functions

- \Rightarrow For episodic functions/env. \rightarrow we can use the start value meaning that starting from s_1 , the value V_1 that we can get.
- ① $J_1(\theta) = V^{\pi_{\theta}}(s_1) = E_{\pi_{\theta}}(G_1) \rightarrow$ total reward
- (If there is a notion of start state).

- ② For continuing env \rightarrow use the average value from that point onwards

$$J_{av}V(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

- ③ or average reward per time step

$$J_{av}R(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(a, s) R_s^a$$

- ④ Same policy gradient \rightarrow works for all. 3.
only distribtz functz changes $\rightarrow d^{\pi_\theta}(s)$.
pn. of being in a state.

Policy-based RL \rightarrow optimization problem.

\rightarrow Find θ to maximize $J(\theta)$

We'll do sequential optimizatz \rightarrow not wait until the very end before optimizing \rightarrow learning within the lifetime of the agent.

Finite difference PG \rightarrow search for a local max. in $J(\theta)$ by ascending the gradient of the policy w.r.t θ .

$$\Delta\theta = \alpha \nabla_\theta J(\theta)$$

$$\nabla_\theta J(\theta) \rightarrow \text{PG}$$

$$= \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

For each dimension k ,
perturb the parameters slightly and take difference \rightarrow

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \mu u_k) - J(\theta)}{\mu}$$

$u_k \rightarrow$ unit vector with k^{th} component = 1
(rest = 0)

\rightarrow n evaluations for n -dimensions

\rightarrow noisy & inefficient

\rightarrow works for arbitrary policy \rightarrow need not be differentiable

(Finite differences).

- FD → gives numerical estimate.
→ collapses for high # dimensions

MC Policy Gradient

- assume policy differentiable whenever it is non-zero
→ Goal: compute gradient analytically
→ assume: know the gradient eg. NN / softmax etc.
as we created it

→ Likelihood ratios → $\nabla_{\theta} \pi_{\theta}(s, a) = \pi_{\theta}(s, a) \frac{\nabla_{\theta} \pi_{\theta}(s, a)}{\pi_{\theta}(s, a)}$

we want to take the expectatⁿ.

$= \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)$

→ score → $\nabla_{\theta} \log \pi_{\theta}(s, a)$ → max likelihood.
functⁿ ↓ we want to follow this gradient.

eg. Softmax policy → smoothly parameterized policy

→ linear combinations of features for actions,
Prob $\propto e^{\phi^T \theta}$ → linear softmax policy
(discrete domain) we choose the one with higher prob.

$\nabla_{\theta} \log \pi_{\theta}(s, a) = \phi(s, a) - E_{\pi_{\theta}}[\phi(s, \cdot)]$
(score functⁿ) actual action feature average feature for all actions we might have taken

Gaussian policy

$\mu(s) = \phi(s, a)^T \theta$

→ fixed / parameterized

Policy is gaussian → $a \sim \mathcal{N}(\mu(s), \sigma^2)$

score functⁿ → $\nabla_{\theta} \log \pi_{\theta}(s, a) = \frac{(a - \mu(s)) \phi(s)}{\sigma^2}$
how much more than usual is being done.

One-step MDP

- starting in state $s \sim d(s)$
- terminate after one time step with reward $r = R_{s,a}$
- use likelihood ratios: (pick action that gets more reward)

$$J(\theta) = E_{\pi_\theta} [r] = \sum_{s \in \mathcal{S}} d(s) \sum_a \pi_\theta(s,a) R_{s,a}$$

(all 3 objective functions become the same here)

$$\nabla_\theta J(\theta) = \sum_s d(s) \sum_{a \in \mathcal{A}} \pi_\theta(s,a) \nabla_\theta \log \pi_\theta(s,a) R_{s,a}$$

to get more reward, just more towards the direction given by score times reward. \rightarrow actual that we saw

take action, compute score \times reward \rightarrow gradient step. \rightarrow sampled reward (model-free)

Now, same in multistep MDP.

- replace immediate reward with value function which gives true policy gradient.

$$\nabla_\theta J(\theta) = E_{\pi_\theta} [\nabla_\theta \log \pi_\theta(s,a) Q^{\pi_\theta}(s,a)]$$

→ adjust such that it does more of the good things & less of the bad things

policy gradient theorem

For supervised learning \rightarrow no value function term (main difference)

Monte Carlo Policy Gradient (REINFORCE)

→ update parameters by SGD

→ sample expectation.

→ we return V_t → unbiased sample of $\mathbb{E}^{\pi_\theta}(s_t, a_t)$

$$\Delta \theta_t = \alpha \nabla_{\theta} [\log \pi_{\theta}(s_t, a_t)] V_t.$$

REINFORCE

Initialize θ arbitrary

for each episode $\{s_1, a_1, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$ do

for $t = 1 \dots T$ do

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) \underbrace{V_t}_{\text{from time step } t \text{ onwards}}$$

end for

end for

return θ

end function

from time step t onwards
(the same we have been using).

MC-PG → curve is smooth
slow 100 million iterations

Actor Critic methods (AC)

MC-PG → still high variance | noisy

$$Q_w(s, a) = \mathbb{E}^{\pi_\theta}(s, a)$$

main idea
instead of using
return, use a
critic for learning
↓
VFA

2 sets of parameters

Critic → updates action-value function parameters w .

Actor → updates policy parameters θ , as judged by critic

↓
doing stuff
in the world.

∴ Approximate PG. (APG)

$$\nabla_{\theta} J(\theta) \approx \mathbb{E}_{\theta} (\nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a))$$

$$\Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)$$

Critic \rightarrow policy evaluation
 \rightarrow How good is current policy π_θ for current param θ .
 • MC • Least squares
 • TD(0)

linear value fn approx $\rightarrow Q_w(s, a) = \phi(s, a)^T w$
 Action value
 Action critic
 Critic \rightarrow update w by TD(0) (linear)
 Actor \rightarrow update θ by policy gradient

function QAC.

Initialize s, θ

sample $a \sim \pi_\theta(s)$

for each step do

sample reward $r = R_s^a$, sample next state s'

sample $a' \sim \pi_\theta(s')$

$$\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$$

$$\theta = \theta + \alpha \nabla_\theta \log \pi_\theta(s, a) \delta Q_w(s, a)$$

$$w = w + \beta \delta \phi(s, a)$$

$$a \leftarrow a', s \leftarrow s'$$

end for

end function

(online algorithm)

initialize $\pi_\theta \rightarrow$ arbitrarily

Bias in actor-critic algorithms

\rightarrow using VFA and approximating the PG \rightarrow introduces bias

\rightarrow But choosing the VFA carefully, we can avoid the bias to get the true PG. \rightarrow compatible

2 conditions:

$$1) \nabla_w Q_w(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

2) w minimizes the mean square error:

$$\epsilon = \mathbb{E}_{\pi_\theta} [Q_w(s, a) - \phi_\pi(s, a)]^2$$

Then PG → exact.

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_{\omega}(s, a)]$$

$$\nabla_{\omega}(E) = 0,$$

$$\mathbb{E}[(Q_{\omega}(s, a) - Q_{\pi_{\theta}}(s, a)) \nabla_{\omega} Q_{\omega}(s, a)] = 0$$

$$\Rightarrow \mathbb{E}_{\pi_{\theta}}[(\quad) \nabla_{\theta} \log \pi_{\theta}(s, a)] = 0 \Rightarrow \text{Proof}$$

$$\Rightarrow \mathbb{E}_{\pi_{\theta}}[Q_{\omega}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)] = \mathbb{E}_{\pi_{\theta}}[Q_{\pi_{\theta}}(s, a) \nabla \log(\cdot)]$$

Softmax policy → get global optimum. in table lookup state.

Tricks

① Reduce variance using Baseline without changing expectation (ie. ascen div²).

Subtract Baseline

$$\mathbb{E}_{\pi_{\theta}}(\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)) = \sum_s d(s) \sum_a \nabla_{\theta} \pi_{\theta}(s, a) B(s)$$

choose such that $\rightarrow \mathbb{E}_{\pi_{\theta}}(\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)) = \sum_s d^{\pi_{\theta}}(s) B(s) \nabla_{\theta} \left(\sum_a \pi_{\theta}(s, a) \right)$

and just reduces variance $\rightarrow \mathbb{E}_{\pi_{\theta}}(\nabla_{\theta} \log \pi_{\theta}(s, a) B(s)) = 0$

\downarrow
 $\text{sum} = 1$
 $\nabla_{\theta} 1 = 0$

so, can add/ subtract any term of this form.

Good baseline $\Rightarrow B(s) = V^{\pi_{\theta}}$ (state value function).

advantage function $A^{\pi_{\theta}}(s, a) = Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s, a)$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\nabla_{\theta} \log \pi_{\theta}(s, a) A^{\pi_{\theta}}(s, a) \right]$$

How much better than usual to take that action.

How to adjust policy to achieve that action.

$A^{\pi_{\theta}}(s, a) \rightarrow$ significantly reduces variance

Ways to estimate:

① estimate both V and $Q \rightarrow$ take difference

$$\textcircled{2} V_V(s) = V^{\pi_{\theta}}(s)$$

$$Q_{\omega}(s, a) = Q^{\pi_{\theta}}(s, a)$$

$$A(s, a) = Q_{\omega} - V_V.$$

update both value functions by TD

② TD error \rightarrow sample of the advantage function

$$\delta^{\pi_\theta} = r + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)$$

$$E_{\pi_\theta}[\delta^{\pi_\theta} | s, a] = E_{\pi_\theta}[r + \gamma V^{\pi_\theta}(s') | s, a] - V^{\pi_\theta}(s)$$

$$= Q^{\pi_\theta}(s, a) - V^{\pi_\theta}(s)$$

TD error is unbiased sample of advantage function (AF)

assuming we have the true values.

$$\nabla_\theta J(\theta) = E_{\pi_\theta}[\nabla_\theta \log \pi_\theta(s, a) \delta^{\pi_\theta}]$$

approximate TD error $\Rightarrow s = r + \gamma V(s') - V(s)$

only requires 1 set of parameters $\rightarrow \underline{V}$

critic at ^{diff.} time scales:

$\rightarrow MC \rightarrow V^t$ is target. $\rightarrow TD(\lambda) \rightarrow \underbrace{V_t^\lambda}_{\text{target}}$

$\rightarrow TD(0) \rightarrow r + \gamma V(s')$

use backward view \rightarrow eligibility traces

$$s_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$e_t = \gamma \lambda e_{t-1} + \phi(s_t)$$

$$\Delta \theta = \alpha s_t e_t$$

Same idea for actors

apply PG with ET.

$$s = r + \gamma V(s_{t+1}) - V(s_t)$$

$$e_{t+1} = \gamma \lambda e_t + \nabla_\theta \log \pi_\theta(s, a)$$

$$\Delta \theta = \alpha s e_{t+1}$$

apply update online

$$\Delta \theta = \alpha [V_t^\lambda - V(s_t)] \nabla_\theta \log \pi_\theta(s, a)$$

introduce bias by bootstrapping to reduce variance

Alternative PG directions

- Gradient-ascent algo: → follow ^{can} any ascent dir
- good dir → significantly speed convergence

so far, stochastic policies → can be very noisy to sample and estimate noise → really hard to estimate.

Natural policy gradient → parameterizes independent.
→ find ascent dir closest to vanilla gradient, when changing policy by a small fixed amount!

$$\nabla_{\theta}^{\text{nat}} \pi_{\theta}(s, a) = \left(G_{\theta}^{-1} \right) \nabla_{\theta} \pi_{\theta}(s, a)$$

work directly on ~~state~~ limiting case, rather than adding noise in policy

G_{θ} → Fisher information matrix

→ start with deterministic policy to get small objective

after taking limiting case, noise → 0.

$$G_{\theta} = E \pi_{\theta} \left[\nabla_{\theta} \log \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a)^T \right]$$

Flashback to compatible value functions.

$$\nabla_{\omega} Q_{\omega}(s, a) = \dot{Q}(s, a) = \nabla_{\theta} \log \pi_{\theta}(s, a) \quad \text{---}$$

⇒ score = features

In case of deterministic policy gradient → directly take expected gradient of Q function.

→ limiting case of stochastic AC.