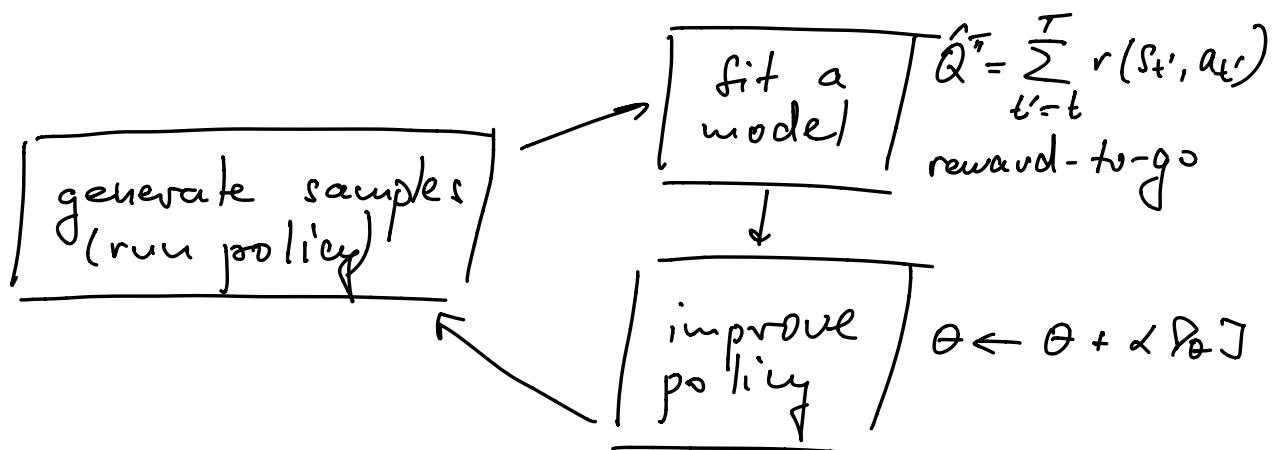


Actor Critic Methods

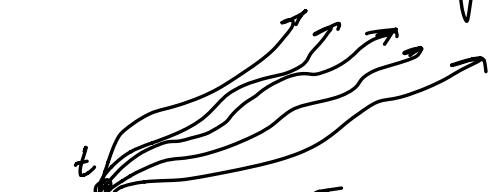
last time : policy gradient (PG)

REINFORCE

1. sample $\{T_j\}_j$ from π_θ
2. estimate $D_\theta J(\theta)$
3. update : $\theta \leftarrow \theta + \alpha D_\theta J(\theta)$



today : discuss better ways to estimate \hat{Q}^π
 → consider trajectory space



$$\hat{Q}_{i,t}^\pi = \sum_{t'=t}^T r(s_{t'}^i, a_{t'}^i) \rightarrow \text{high variance in estimating } D_\theta J$$

$$\text{want : } Q_t^\pi \approx \sum_{t'=t}^T E_{\pi_\theta} [r(s_{t'}, a_{t'})] \text{ expected reward-to-go}$$

- smaller variance
- more consistent gradient estimates
- learn faster
- baseline: does not change \hat{P}_θ]
 - use state-dependent baseline
- $b = V(s_t) = \mathbb{E}_{a_t \sim \pi_\theta} [Q(s_t, a_t)]$
leaves \hat{P}_θ invariant (HW)
- but: expectation over action further reduces variance!

note: action-dep. baselines change \hat{P}_θ]

recall :

$$Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_t, a_t') | s_t, a_t]$$

expected return following policy π_θ
given (s_t, a_t)

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi_\theta} [Q^\pi(s_t, a_t) | s_t]$$

expected return following policy π_θ
given s_t

def : advantage function:

$$A^\pi(s_t, a_t) := Q^\pi(s_t, a_t) - V^\pi(s_t)$$

how much better action a_t is than
the average action in state s_t under π_θ

- using $b = V^\pi(s_t)$, PG update reads

$$\nabla J(\theta) \approx \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \text{Policy}_{\pi_\theta}(a_t^j | s_t^j) \left[\underbrace{\left(\sum_{t'=t}^T r(s_t^{j'}, a_t^{j'}) \right)}_{= A^{\pi_\theta}(s_t^j, a_t^j)} - b \right]$$

$$= \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T \text{Policy}_{\pi_\theta}(a_t^j | s_t^j) A^{\pi_\theta}(s_t^j, a_t^j)$$

- idea: the better we can estimate A^{π_θ} ,
the lower the variance of ∇J

- Value function fitting:

- which one of A^π , Q^π , V^π should we fit?

$$\begin{aligned} Q^\pi(s_t, a_t) &= r(s_t, a_t) + \sum_{t'=t+1}^T \mathbb{E}_{\pi_\theta}[r(s_{t'}, a_{t'}) | s_t, a_t] \\ &= r(s_t, a_t) + \mathbb{E}_{s_{t+1} \sim p(s_{t+1} | s_t, a_t)} [V^\pi(s_{t+1})] \end{aligned}$$

use single sample estimate
for p

$$\approx r(s_t, a_t) + V^\pi(s_{t+1})$$

$$\Rightarrow A^\pi(s_t, a_t) \approx r(s_t, a_t) + V^\pi(s_{t+1}) - V^\pi(s_t)$$

\rightarrow only V^π occurs! \rightarrow fit V^π

use a DNN with parameters γ

$$s \rightarrow \text{DNN}_\gamma \rightarrow V_\gamma^\pi(s)$$

- policy evaluation: what do we fit V_γ^π to?

$$V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_t]$$

$$\mathcal{J}(\theta) = \mathbb{E}_{s_0 \sim p(s_0)} [V^\pi(s_0)]$$

1) MC policy evaluation (same as PG, last time)

$$V^\pi(s_t) \approx \sum_{t'=t}^T r(s_{t'}, a_{t'})$$

single sample estimate
(unlikely to repeat trajectory
twice)

why meaningful? s

$$V_\gamma^\pi(s) \approx V_\gamma^\pi(s)$$

different trials visit
different states but
NN can interpolate/
generalize

training data set for V_p^π :

$$\{ s_t^j, y_t^j = \sum_{t'=t}^T r(s_{t'}^j, a_{t'}^j) \}_{j=1}^N$$

↑ ↑
data points "labels"

→ use mean-square loss (similar to regression)

$$\text{loss: } l(\varphi) = \frac{1}{2} \sum_{j=1}^N \| V_p^\pi(s_j) - y_j \|^2$$

→ unbiased estimator ✓

→ single sample estimates → large variance

2) bootstrap estimate (TD learning estimate)

$$V^\pi(s_t^j) \approx r(s_t^j, a_t^j) + \underbrace{\sum_{t'=t+1}^T \mathbb{E}_{\pi_\theta} [r(s_{t'}, a_{t'}) | s_{t+1}^j]}_{\text{neglect } \mathbb{E}_p(s_{t+1} | s_t, a_t)}$$

$$\mathbb{E}_p(s_{t+1} | s_t, a_t) \approx \underbrace{V^\pi(s_{t+1}^j)}_{\text{don't know this value yet at step } t} \approx \underbrace{V_p^\pi(s_{t+1})}_{\substack{\text{plug in previous iteration estimate} \\ \parallel}}$$

biasd estimator

but lower variance

training data:

$$\{s_t^j, y_t^j = r(s_t^j, a_t^j) + V_\phi^\pi(s_{t+1}^j)\}_{j=1}^N$$

→ same loss as before $\rightarrow \pi_\theta \rightarrow V_\phi^\pi$

- Algorithm: Offline Actor Critic method

(PG with value function estimation)

1. sample $\{s_t^j, a_t^j\}$ from $\pi_\theta \rightarrow$ go until end of episode!

2. fit $V_\phi^\pi(s)$ to the sampled data:

$$L(\phi) = \frac{1}{2} \sum_j \|V_\phi^\pi(s_j) - y_j\|^2$$

a) MC estimate: $y_t^j = \sum_{t=t'}^T r(s_{t'}, a_{t'})$

OR

b) bootstrap estimate: $y_t^j = r(s_t^j, a_t^j) + \gamma V_\phi^\pi(s_{t+1}^j)$

3. evaluate advantage function:

$$A^\pi(s_t^j, a_t^j) \approx r(s_t^j, a_t^j) + \gamma V_\phi^\pi(s_{t+1}^j) - V_\phi^\pi(s_t^j)$$

4. compute the $D_\theta J$:

$$D_\theta J \approx \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T D_\theta \log \pi_\theta(a_t^j | s_t^j) \cdot A^\pi(s_t^j, a_t^j)$$

$$5. \theta \leftarrow \theta + \alpha \nabla \mathcal{J}(\theta)$$

- problem: keep adding values to $V^\pi(s)$
 $\rightarrow V^\pi(s)$ could grow indefinitely

\rightarrow use a discount factor $\gamma \in [0, 1]$

$$y_t^i \approx r(s_t^i, a_t^i) + \gamma V_\phi^\pi(s_{t+1}^i)$$

$$A^\pi(s_t^i, a_t^i) = r(s_t^i, a_t^i) + \gamma V_\phi^\pi(s_{t+1}^i) - V_\phi^\pi(s_t^i)$$

γ decreases reward-to-go
 \Rightarrow reduces variance

NB: in AC methods we need $\gamma < 1$
 $(\gamma \approx 0.99)$ even in episodic tasks!

intuition: rewards in the immediate future count more than rewards in distant future

observation: bootstrap estimate does not require to finish episode in order to be computed

- Algorithm : Online Actor-Critic method

1. take action $a \sim \pi_\theta(a|s)$, get (s, a, r, s')

2. update V_φ^π using bootstrap/TD target
 $y(s) = r + \gamma V_\varphi^\pi(s')$: $\ell(\varphi) = \frac{1}{2} \|V_\varphi^\pi(s) - y(s)\|^2$

3. estimate advantage function

$$A^\pi(s, a) \approx r(s, a) + \gamma V_\varphi^\pi(s') - V_\varphi^\pi(s)$$

4. $D_\theta J \approx D_\theta \log \pi_\theta(a|s) A^\pi(s, a)$ no sums
over trajectories
or time steps!

5. $\theta \leftarrow \theta + \alpha D_\theta J$

- comparison AC vs. PG :

• AC :

$$D_\theta J = \frac{1}{N} \sum_j \sum_t D_\theta \log \pi_\theta(a_t^j | s_t^j) \times \\ \times [r(s_t^j, a_t^j) + \gamma V_\varphi^\pi(s_{t+1}^j) - V_\varphi^\pi(s_t^j)]$$

→ lower variance (due to critic V_φ^π)

→ biased (if critic is not perfect/exact)

- PG

$$D_\theta] \approx \frac{1}{N} \sum_j \sum_t \text{Policy}(\hat{\pi}_\theta(a_t^j | s_t^j)) \times \\ \times \left[\left\{ \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}^j, a_{t'}^j) \right\} - b \right]$$

→ no bias (MC target)

→ high variance (single-sample estimate)

- $b = V_p^\pi(s)$: state-dep. estimate

→ no bias

→ lower variance (not as much as AC)

- practical implementation:

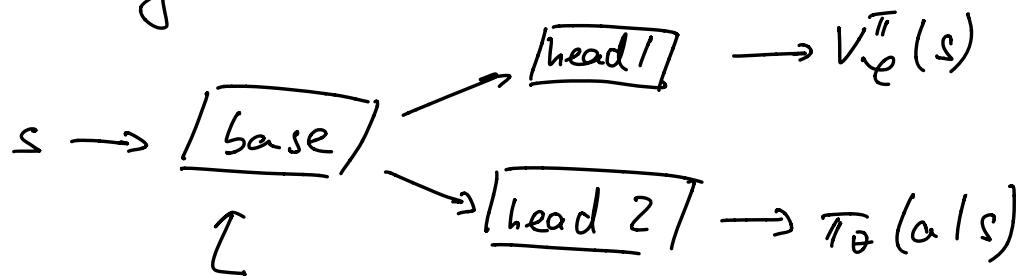
- need to parametrize both π_θ & V_p^π

- i) network architecture

- i) two independent deep neural nets

$$\begin{aligned} \text{critic: } s \rightarrow DNN_p \rightarrow V_p^\pi(s) \\ \text{actor: } s \rightarrow DNN_\theta \rightarrow \pi_\theta(a|s) \end{aligned} \quad \left. \begin{array}{l} \text{stable} \\ \text{training } \checkmark \\ \text{no shared} \\ \text{features } \times \end{array} \right\}$$

ii) single base + two heads



- shared features are possible ✓
- shared gradients in base layer can be of different scale for V_ϕ^π , π_θ
→ fine-tune learning rates

2) issue with online AC:

uses a minibatch size of 1!

→ use parallelization (JAX → pmap, vmap)

a) synchronized:

parallel agents, share π_θ , V_ϕ^π

$\downarrow \downarrow \downarrow \downarrow$ barrier, get #agents

batch $\{(s, a, r, s')\}_{j=1}^J$

update parameters Θ, φ

$\downarrow \downarrow \downarrow \downarrow$ barrier → new minibatch, etc.

b) asynchronous parallel AC
→ Berkeley video lectures