# Deep Reinforcement Learning

- so far : discussed <u>action-value</u> methods for RL
  $\rightarrow$ policy defined from $V_\pi(s)$ or $q_\pi(s,a)$

$V_\pi(s) \xrightarrow{\text{DRL}} V_\theta(s)$   approximate

$q_\pi(s,a) \longrightarrow q_\theta(s,a)$   value functions &

$\pi(a/s) \longrightarrow \pi_\theta(s,a)$   policy using a

ML model with

params $\theta$: weights

& biases

## Policy Gradient Methods

- algorithms that learn a parametrized
  policy directly : $\pi_\theta(s|a)$

• value function may be present, but it
  is <u>not</u> required to select actions
   ( $\rightarrow$ actor - critic algorithms)

<u>idea</u> : parametrize policy by some variational
  parameters $\theta \in \mathbb{R}^d$
  $$\pi(a/s) \approx \pi_\theta(a/s)$$

1) <u>which</u> <u>variational ausatz</u>/model do we
<u>use</u> for $\theta \longmapsto \pi_\theta(a/s)$ ?

e.g. i) if action space $A$ is <u>discrete</u>:
<u>softmax policy</u>:

$$\pi_\theta(a/s) := \frac{e^{\beta h_\theta(s,a)}}{\sum_{a \in A} e^{\beta h_\theta(s,a)}}$$

$\beta$ : (inverse) temperature
  • $\beta \to \infty$ : policy is greedy
  • $\beta \to 0$ : equiprobable policy

$h_\theta(s,a)$ : some parametrized function
    on $S \times A$, e.g. DNN or CNN
          or anything suitable
        to problem of study

e.g. Gridworlds, Atari games, board games

ii) <u>continuous</u> action space $A$

e.g. Gaussian policy:

$$\pi_\theta(a/s) := \frac{1}{\sqrt{2\pi \sigma_\theta^2(s)}} e^{-\frac{(a - \mu_\theta(s))^2}{2 \sigma_\theta^2(s)}}$$

-2-

$\theta \longmapsto \mu_\theta$   parametrize the mean $\mu_\theta$
$\theta \longmapsto \sigma_\theta$   & variance $\sigma_\theta^2$  of the
Gaussian policy

e.g.   $s \longrightarrow DNN_\theta \Big\langle \begin{array}{c} \nearrow \mu_\theta \\ \searrow \sigma_\theta \end{array} \Big\} \rightarrow \overline{\pi_\theta}$

issue: <u>unimodality</u>:

opt. $\pi_*$
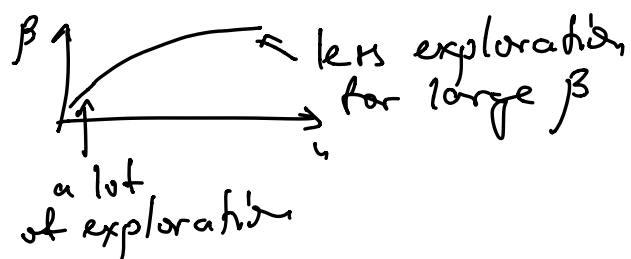
$\rightarrow$ single Gaussian insuff.

$\rightarrow$ Gaussian mixture models, correlated Gaussian
etc.

— <u>possible</u> <u>advantages</u> of policy gradient
methods (over value-function methods):

$\rightarrow \pi_\theta(a|s)$ can approach a deterministic policy
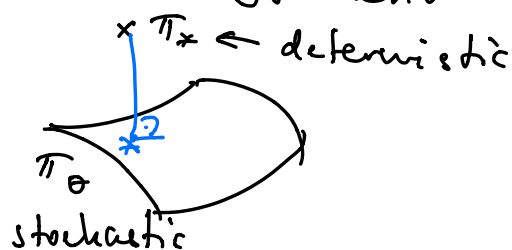(not possible if we use $\varepsilon$-greedy policies,
e.g. SARSA, Q-learning etc.)

• use (inv.) temp parameter $\beta$ to control
stochasticity using a schedule:

$\beta = \beta(n)$ , $n$: episode number

β ← less exploration for large β

a lot of exploration

- using a schedule we can enforce convergence of the algorithm (but to what?!)

→ parametrization enables (in principle) an arbitrary policy (can have a stochastic optimal policy)
  e.g. reward is uncertain
       or env. is stochastic



× $\pi_*$ ← determistic

$\pi_\theta$
stochastic

→ learning the policy might just be easier than learning a value function!

→ can use policy to build in prior knowledge about problem (pretraining, etc.)

2) how do we find optimal parameters $\theta$?
   s.t. $\pi_\theta(a|s) \approx \pi_*(a|s)$ for a given task?

   $\rightarrow$ assume that $\theta \longmapsto \pi_\theta(a|s)$ is diff'ble

   $\Rightarrow$ find those params $\theta$ which minimize
   expected returns $\mathbb{E}_{\pi_\theta}[G_t | S_t = s_0]$

   e.g. use gradient descent/ascent

   ## Policy Gradient Theorem

   - RL objective : total expected returns
     starting from state $s_0$
     & following policy $\pi_\theta$ afterwards

   $$J(\theta) := \mathbb{E}_{\pi_\theta}[G_t | S_t = s_0]$$

   $$= V_{\pi_\theta}(s_0) \quad \text{value function of } s_0$$

   — compute gradient of the RL objective
   wrt. varl. params. $\theta$

$$\nabla_\theta V_{\pi_\theta}(s) = \nabla_\theta \sum_a \pi_\theta(a|s) \, q_{\pi_\theta}(s,a)$$

$$= \sum_a \left( \nabla_\theta \pi_\theta(a|s) \right) q_{\pi_\theta}(s,a) + \pi_\theta(a|s) \nabla_\theta \, q_{\pi_\theta}(s,a)$$

$$= \sum_a \Big\{ \left[ \nabla_\theta \pi(a|s) \right] q_{\pi_\theta}(s,a) +$$

$$+ \pi_\theta(a|s) \nabla_\theta \sum_{s',r} p(s',r|s,a) \left[ r + V_{\pi_\theta}(s') \right] \Big\}$$

$$\underset{\sum_r p(s',r|s,a) = p(s'|s,a)}{} \qquad \underset{\text{indep. of } \theta}{}$$

$$= \sum_a \Big\{ \left[ \nabla_\theta \pi(a|s) \right] q_{\pi_\theta}(s,a) +$$

$$+ \pi_\theta(a|s) \sum_{s'} p(s'|s,a) \, \nabla_\theta V_{\pi_\theta}(s') \Big\}$$

unroll

$$\downarrow$$

$$= \sum_a \Big\{ \left[ \nabla_\theta \pi_\theta(a|s) \right] q_{\pi_\theta}(s,a) +$$

$$\pi_\theta(a|s) \sum_{s'} p(s'|s,a) \sum_{a'} \Big[ \left( \nabla_\theta \pi_\theta(a'|s') \right) q_{\pi_\theta}(s',a') +$$

$$+ \pi_\theta(a'|s') \sum_{s''} p(s''|s',a') \, \nabla_\theta V_{\pi_\theta}(s'') \Big] \Big\}$$

$$= \ldots =$$

$$= \sum_{x \in S} \sum_{k=0}^{\infty} Pr_{\pi_\theta}(s \to x; k) \sum_a \left[ \nabla_\theta \pi_\theta(a|x) \right] q_{\pi_\theta}(x,a)$$

where $Pr_{\pi_\theta}(s \to x; k)$ is the prob. to transition

from $s \to x$ in $k$ steps after following $\pi_\theta$

$\to$ evaluate on initial state $s_0$:

$$D_\theta \, J(\theta) = \sum_{s \in S} \underbrace{\left( \sum_{k=0}^{\infty} P_{v_{\pi_\theta}} (s_0 \to s; k) \right)}_{=: \, \gamma_{\pi_\theta}(s)} \sum_a \left( D_\theta \pi_\theta(a|s) \right) q_{\pi_\theta}(s,a)$$

issue: $\gamma_{\pi_\theta}(s)$ is $\underline{not}$ a prob. distr.
(not normalized)

$$D_\theta \, J(\theta) = \underbrace{\sum_{s'} \gamma_{\pi_\theta}(s')}_{\substack{= \text{const.} \\ \text{(indep. of } \theta)}} \sum_s \underbrace{\frac{\gamma_{\pi_\theta}(s)}{\sum_{s''} \gamma_{\pi_\theta}(s'')}}_{=: \, \mu_{\pi_\theta}(s) \text{ normalized!}} \sum_a D_\theta \pi_\theta(a|s) \times q_{\pi_\theta}(s,a)$$

$$D_\theta \, J(\theta) \propto \sum_s \mu_{\pi_\theta}(s) \sum_a D_\theta \pi_\theta(a|s) \times q_{\pi_\theta}(s,a)$$

overall const in GD is irrevaut in practice
b/c it can be absorbed in the learning rate/
step size of GD.

$$\boxed{D_\theta \, J(\theta) \propto \sum_s \mu_{\theta}(s) \sum_a D_\theta \pi_\theta(a|s) \times q_{\pi_\theta}(s,a)}$$

— $\underline{problem}$: $\mu_{\pi_\theta}(s)$ requires knowledge of
the transition prob. $p(s'|s,a)$
$\to$ it requires a model

→ way out : MC methods !

need to write $\nabla_\theta J(\theta)$ as an expectation
over trajectories $\tau = (s_0, A_0, S_1, A_1, \ldots)$

$$\nabla_\theta J = \sum_{s,a} \mu_{\pi_\theta}(s) \; q_{\pi_\theta}(s,a) \; \pi_\theta(a/s) \; \underbrace{\frac{\nabla_\theta \pi_\theta(a/s)}{\pi_\theta(a/s)}}_{= \nabla_\theta \log \pi_\theta(a/s)}$$

$$= \sum_{s,a} \underbrace{\mu_{\pi_\theta}(s) \pi_\theta(a/s)}_{\substack{\text{probability for} \\ \text{pair } (s,a) \text{ to occur} \\ \text{under } \pi_\theta \text{ \& } P}} \nabla_\theta \log \pi_\theta(a/s) \times q_{\pi_\theta}(s,a)$$

—def : prob. of a trajectory $\tau$ to occur
under policy $\pi_\theta$ \& trans. prob. $p(s'/s,a)$

$$\underline{P}(\tau) := p(s_0) \pi_\theta(a_0/s_0) p(s_1/s_0,a_0) \pi_\theta(a_1/s_1) p(s_2/s_1,a_1) \ldots$$

$$= p(s_0) \prod_{t=0}^{T-1} \pi_\theta(a_t/s_t) p(s_{t+1}/s_t, a_t)$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \underline{P}} \left[ \nabla_\theta \log \pi(\tau) \times G(\tau) \right]$$

where $\pi_\theta(\tau) := \prod_{t=0}^{T-1} \pi_\theta(a_t/s_t) \; ; \; G(\tau) = \sum_{t=0}^{T-1} r(s_t, a_t)$

$$\tau = (s_0, a_0, s_1, a_1, \ldots, a_{T-1}, s_T)$$

-8-

— generate/sample $N$ trajectories using MC from policy $\pi_\theta$

$$D_\theta J(\theta) \approx \frac{1}{N} \sum_{j=1}^{N} D_\theta \log \pi_\theta(\tau^j) \times G(\tau^j) \quad \substack{\text{trajectory} \\ \text{sampled using} \\ \text{MC}}$$

$$= \frac{1}{N} \sum_{j=1}^{N} \left[ \sum_{t=1}^{T} D_\theta \log \pi_\theta(a_t^j / s_t^j) \right] \left[ \sum_{t'=1}^{T} r(s_t^j, a_t^j) \right]$$

$\underset{\text{within the episode}}{\text{time step}}$

— update rule for parameters $\theta$ using SGA:

$$\theta \longleftarrow \theta + \alpha \frac{D_\theta J(\theta)}{\uparrow}$$

ascent (maximize return)      learning rate/step size      use MC estimate

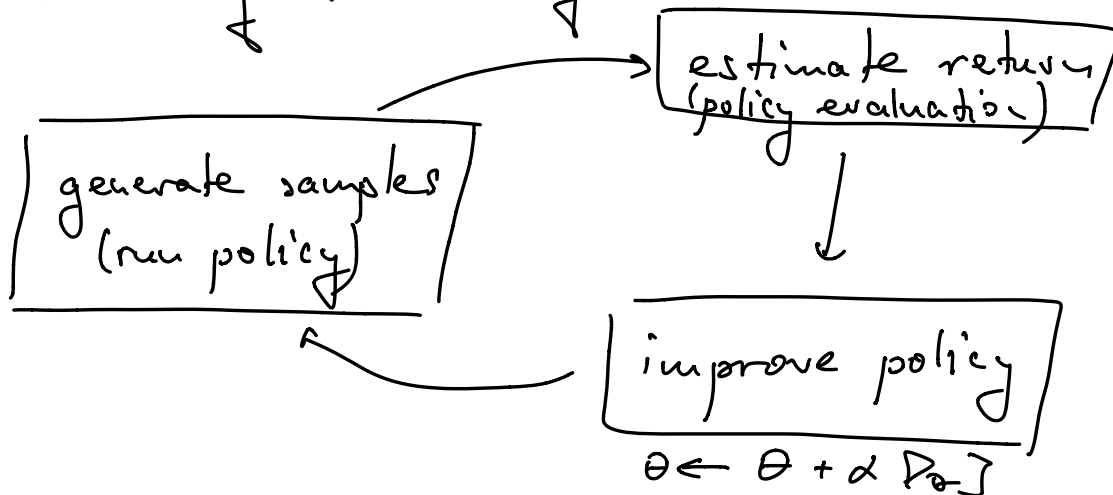— alternative derivation in Berkeley lectures

— <u>intuition about PG</u>:

→ PG : "gradient ascent for RL"

→ PG makes higher-reward trajectories more likely

→ PG : trial & error learning ( ⇒ MC)
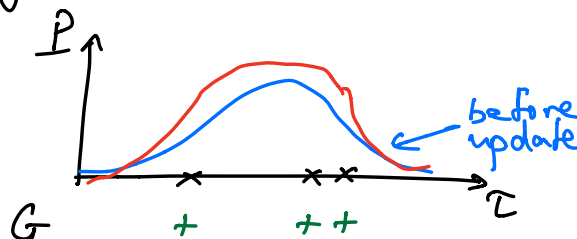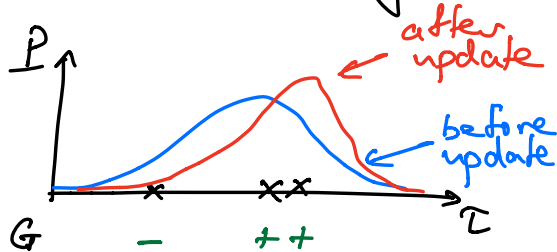
— Anatomy of RL algorithms



$$\theta \leftarrow \theta + \alpha \, \nabla_\theta J$$

— <u>REINFORCE</u> algorithm for PG :

1. sample $\{\tau^i\}_{i=1}^N$ using MC from current policy $\pi_\theta$ (run policy)

2. estimate $\nabla_\theta J(\theta)$

3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$ (improve policy)

— What can go <u>wrong</u> with REINFORCE?



$$r \rightarrow r + \text{large const.}$$

$\rightarrow$ hints at a high variance problem

— how do we <u>reduce the variance</u> in PG?

1) take into account <u>causality</u>:

$$\nabla_\theta J \approx \frac{1}{N} \sum_{j=1}^{N} \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta (a_t^j | s_t^j) \sum_{t'=1}^{T} r(s_{t'}^j, a_{t'}^j)$$

past actions affect future states & rewards
but not vice-versa!

$\pi_{t'}$ cannot affect $r_t$ at $t < t'$

$$\sum_{t'=1}^{T} r(s_{t'}^j, a_{t'}^j) \longrightarrow \sum_{t'=t}^{T} r(s_{t'}^j, a_{t'}^j) = q_t^j$$

reduces variance b/c we have to sum
ups fewer members

$q_t^j$ : <u>reward-to-go</u>

2) <u>baselines</u> :

<u>recall</u> : i) PG makes higher reward trajectories
more probable (intuition)

ii) e.g. $r \longrightarrow r + 10^9 \Rightarrow$ all trajectories
seem to be good

want : make traj. which are better
than average more likely $\&$
trajectories worse than average
less likely.

$$b := \frac{1}{N} \sum_{\hat{f}=1}^{N} r(\tau\hat{f}) \quad \text{average return}$$

want:

$$\nabla_\theta J \approx \frac{1}{N} \sum_{\hat{f}} \nabla_\theta \log \pi_\theta(\tau\hat{f}) [G(\tau\hat{f}) - b]$$

issue:: that's not the original gradient!

look at:

$$\mathbb{E}_{\tau \sim P} [b \nabla_\theta \log \pi_\theta(\tau)] = b \mathbb{E}_{\tau \sim P} [\nabla_\theta \log \pi_\theta(\tau)]$$

$$= b \mathbb{E}_{\tau \sim P} \left[ \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} \right]$$

$$= b \int d\tau \, \pi_\theta(\tau) \frac{\nabla_\theta \pi_\theta(\tau)}{\pi_\theta(\tau)} = b \int d\tau \, \nabla_\theta \pi_\theta(\tau)$$

$$= b \nabla_\theta \int d\tau \, \pi(\tau) = b \nabla_\theta \underbrace{\mathbb{E}_{\tau \sim P}[1]}_{=1} = 0$$

$\rightarrow$ adding a baseline does not change
the expectation of the gradient! ($\rightarrow$ baselines
are unbiased)

but it changes its variance

- improved PG with baseline: (lower variance)

$$\nabla_\theta J \approx \frac{1}{N} \sum_{j=1}^{N} \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta (a_t^j \mid s_t^j) \sum_{t'=t}^{T} \left( r(s_{t'}^j, a_{t'}^j) - b \right)$$

— optimal baseline for minimal variance:

recall :

$$\mathrm{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

$$\mathrm{Var}(b) = \mathbb{E}_{\tau \sim P} \left[ (\nabla_\theta \log \pi_\theta (G-b))^2 \right] - \underbrace{\mathbb{E}_{\tau \sim P} \left[ \nabla_\theta \log \pi_\theta (G-b) \right]^2}_{= \mathbb{E}_{\tau \sim P} \left[ \nabla_\theta \log \pi_\theta \, G \right]^2}$$

$$\to \text{indep. of } b\,!$$

$$\frac{\partial}{\partial b} \mathrm{Var}(b) = \mathbb{E}_{\tau \sim P} \left[ (\nabla_\theta \log \pi_\theta)^2 (-2G + 2b) \right]$$

$$= -2 \mathbb{E}_{\tau \sim P} \left[ (\nabla_\theta \log \pi_\theta)^2 \, G \right] + 2b \, \mathbb{E}_{\tau \sim P} \left[ (\nabla_\theta \log \pi_\theta)^2 \right]$$

$$\frac{\partial}{\partial b} \mathrm{Var}(b) \overset{!}{=} 0$$

$$\Rightarrow b_{\text{optimal}} = \frac{\mathbb{E}_{\tau \sim P} \left[ (\nabla_\theta \log \pi_\theta)^2 \, G \right]}{\mathbb{E}_{\tau \sim P} \left[ (\nabla_\theta \log \pi_\theta)^2 \right]} \quad \propto \text{understand division element-wise}$$

$\to b_{\text{optimal}}$ knows about different directions in $J(\theta)$-landscape

→ <u>drawback</u> : not really used in practice
b/c tedious to implement

— <u>Practical implementation</u> of <u>PG</u> :

→ how do we efficiently compute $\nabla_\theta J$ ?

$$\nabla_\theta J \approx \frac{1}{N} \sum_j \left[ \nabla_\theta \log \pi_\theta (\tau^j) \right] \left( G(\tau^j) - b \right)$$

<u>but</u> : $J \approx \frac{1}{N} \sum_j G(\tau^j)$ indep. of $\theta$

→ cannot use that
as a cost function?

use a <u>pseudo cost fn</u> :
→ a cost/loss function , s.t.

$$\nabla_\theta J_{pseudo}(\theta) = \nabla_\theta J_{exact}(\theta)$$

but $J_{pseudo}(\theta) \neq J_{exact}(\theta)$

$$J_{pseudo}(\theta) \approx \frac{1}{N} \sum_j \log \pi_\theta (\tau^j) \left[ G(\tau^j) - b) \right]$$

— a few <u>implementation</u> <u>tips</u> :
1) use larger batches of data
as compared to supervised learning
( need to reduce variance of $\nabla_\theta J$ )

2) learning rate of ADAM/SGD may
need fine-tuning