



1 MDP. Bellman eq. & operator  $T^\pi V(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} \mathbb{E}_{s' \sim P(\cdot|s,a)} (r(s,a,s') + \gamma V(s'))$

$T^\pi$  is mono., contraction wrt.  $\|\cdot\|_\infty$ .  $T^\pi V^\pi = V^\pi$ .  $V^\pi(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} Q^\pi(s,a)$

(action-value version)  $\rightarrow$  fix point  $V^\pi$ .  $m+1$  form  $T^\pi V = \gamma^\pi + P^\pi V$

$\mathbb{E}^\pi Q(s,a) = \mathbb{E}_{s' \sim P(\cdot|s,a)} (r(s,a,s') + \gamma \mathbb{E}_{a' \sim \pi(\cdot|s')} Q(s',a'))$  and  $(I - \gamma P^\pi)^{-1} \geq I$ .  
 $\mathbb{E}^\pi Q^\pi = Q^\pi$ .  $= I + \gamma P^\pi + \gamma^2 P^{\pi^2} + \dots$

Optim  $V^*(s) \triangleq \max_\pi V^\pi(s)$ .  $\square 1$  (Policy improvement)  $\pi'(a|s) = \begin{cases} 1 & a = \arg \max_a Q^\pi(s,a) \\ 0 & \text{ov.} \end{cases}$ ,  $V^{\pi'} \geq V^\pi$ .  
 $\square 2$  (Optimal existence)

(not unique)  $\exists \pi^* \triangleq \{ \cdot | V^* \text{ s.t. } V^* = V^{\pi^*}$ . Pf:  $V^\pi = T^\pi V^\pi \leq T^{\pi'} V^\pi$ ,  $V^\pi \leq \lim_k (T^{\pi'})^k V^\pi = V^{\pi'}$  (by  $\square 2$ )

Pf:  $V^\pi \leq T^{\pi^*} V^*$ ,  $V^\pi \Rightarrow V^* \leq V^{\pi^*}$ . Bellman optimality eq. & operator  $T V^* = V^*$

$T$  is mono., contraction wrt.  $\|\cdot\|_\infty$ .  $T V(s) = \max_a \mathbb{E}_{s' \sim P(\cdot|s,a)} (r(s,a,s') + \gamma V(s'))$

(a-v version)  $Q^*(s,a) \triangleq \max_\pi Q^\pi(s,a)$ ,  $V^*(s) = \max_a Q^*(s,a)$ .  $\mathbb{E} Q^* = Q^*$ .  
 $= \mathbb{E}_{s \sim P(\cdot|s,a)} (r(s,a,s') + \gamma V^*(s')) = Q^{\pi^*}(s,a)$ .  $\mathbb{E} Q(s,a) = \mathbb{E}_{s' \sim P(\cdot|s,a)} (r(s,a,s') + \gamma \max_a Q(s',a'))$

VI.  $V^{k+1} = T V^k = T^{\pi_{k+1}} V^k$ .  $\|V^k - V^*\|_\infty \leq \gamma^k \|V^0 - V^*\|_\infty$ .  $O(|S||A|)$  per iter.  $\max_a Q(s',a')$   
async. (Gauss-Seidel)  $V(s) \leftarrow T V(s)$ ,  $s=1,2,\dots$

PI.  $\pi_{k+1} = \{ \cdot | Q^{\pi_k} \}$ .  $(T^{\pi_{k+1}} V^{\pi_k} = T V^{\pi_k} \geq V^{\pi_k}$ ,  $\|\cdot\|_\infty$ .  $F(V) := \max_\pi \{ \gamma^\pi + P^\pi V \} - V = 0$ ,  
(first eval  $V^{\pi_k}$ ) by  $\square 1$   $V^{\pi_{k+1}} \geq V^{\pi_k}$ . which is  $\square 1$  Pf) (Pf:  $V^{\pi_{k+1}} \geq T V^{\pi_k} \geq V^{\pi_k}$ ,  $V^{\pi_k} - (\gamma P^{\pi_{k+1}} - I)^{-1} F(V^{\pi_k})$   
is  $\pi_{k+1}$  for  $V^{\pi_k}$ )

full exploration. if use a-v version can be model-free.  $V^* - V^{\pi_k} \leq T^k (V^* - V^0)$ . Jacobian  $= V^{\pi_{k+1}}$

$\circ$  TPI ( $T^{\pi_k} V^{\pi_{k-1}} = V^k$ ,  $\pi_{k+1} = \{ \cdot | V^k \}$ . API.  $\|V^k - V^*\|_\infty \leq \delta$  eval,  
 $m_k \rightarrow \infty$  is PI,  $m_k = 1$  is VI.  $\| \gamma^{\pi_{k+1}} + \gamma P^{\pi_{k+1}} V^k - T V^{\pi_{k+1}} \|_\infty \leq \epsilon$

$\square 3$  (error amplification)  $\pi = \{ \cdot | V$ ,  $V^\pi \geq V^* - \frac{2\gamma}{1-\gamma} \|V - V^*\|_\infty$  ( $\epsilon$ -optim.) Pf:  $\geq$

Strongly polynomial PI, VI, LPI.  $LPI$ .  $\sim \frac{2}{1-\gamma} \|Q - Q^*\|_\infty T$

$\square 4$  (strict progress)  $\{ \pi_k \}$  from PI,  $\exists s$  s.t.  $V^k \geq \frac{1}{1-\gamma} \log \frac{1}{1-\gamma}$ ,  $\pi_k(s) \neq \pi_0(s)$ . Pf:  $\geq$ .  
 $\Rightarrow$  at most  $O(\frac{|S||A|-1}{1-\gamma} \log \frac{1}{1-\gamma})$  iters to optim.



2

MC. variation ctrl  $E f(x) = E(f(x) + c(g(x) - E g(x)))$ , importance sampling  $E f(x) = E \frac{f(y)p(y)}{q(y)}$   
 $c^* = -\frac{Cov(f, g)}{V(g)}$  bias-var tradeoff  $q^*(y) = \frac{f(y)p(y)}{E_x |f| p(x)}$

Robbins-Monro algo.  $x_{k+1} = x_k + \alpha_k (h(x_k) - x_k + e_k)$  for  $h(x) = x$ .  
eg.  $h(x) = \mu$ , SGD  $= x - E_g \nabla f(x; \xi)$  (RM convergence)  $|h'(x)| \leq \gamma < 1, \forall x. \sum \alpha_k = \infty, \sum \alpha_k^2 < \infty$ .  
Q5  $E(e_k | \mathcal{F}_k) = 0, E(e_k^2 | \mathcal{F}_k) < \infty. x_k \xrightarrow{a.s.} x^*$

Q5' (conv. rate) if  $\alpha_k = \frac{1}{k}$ ,  $E \|x_k - x^*\|_2^2 \approx \frac{1}{k} \frac{Cov(f, g)}{V(g)}$  Pf: Dvoretzky Thm. (u.a.s.)  
sto. process  $w_{k+1} = (1 - \alpha_k)w_k + \beta e_k, e.d. \beta \sim \infty, w_k \xrightarrow{a.s.} 0$ . let  $x_k - x^* = w_k$ .

(eval by MC) primitive MC.  $\nabla \tau \sim \pi_k$  from  $(s, a)$  (Pf of Dvoretzky Thm.  $v_k := w_k^2, u_k := \alpha_k^2 w_k^2 + \beta^2 E(e_k^2 | \mathcal{F}_k)$ .  
 $E(v_k | \mathcal{F}_k) \leq (1 - \alpha_k)v_k + u_k$  with  $\sum \alpha_k = \infty, \sum u_k < \infty$ .  
 $S_k := v_k + \sum_{j=1}^k u_j$  is supermg. Doob  $\Rightarrow S_k \xrightarrow{a.s.} \infty$ .  
 $E(S_{k+1} | \mathcal{F}_k) \leq S_k - \alpha_k u_k \Rightarrow \sum \alpha_k v_k \leq E S_k < \infty$ .

(policy impr by a-v) learn  $Q^k(s, a) = \frac{1}{n} \sum_{i=1}^n \gamma^i r_t^i$   
 $\pi_{k+1} = \arg \max_a Q^k(s, a)$

use subtraj. first/every visit.  
every visit MC.  $\nabla \tau \sim \pi$  for  $t = T-1 \rightarrow 0$   
eval  $G = 0$   
(or in incremental way  $N(s_t) += 1, G(s_t) += G, V(s_t) = G(s_t) / N(s_t)$ )  
 $V(s_t) += \frac{G - V(s_t)}{N(s_t)}$  end

MC learn with greedy  $\nabla \tau \sim \pi_k$  for  $t = T-1 \rightarrow 0$   
(may fail since lack of exploration) if  $(s_t, a_t)$  not in  $(s_0, a_0, \dots, a_{t-1})$   
 $N(s_t, a_t) += 1; Q(s_t, a_t) += \frac{G - Q(s_t, a_t)}{N(s_t, a_t)}$   
 $\pi_{k+1}(a|s_t) = \frac{Q(s_t, a)}{\sum_{a'} Q(s_t, a')}$  end

$\epsilon$ -greedy.  $\pi'(a|s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A|} & a = \arg \max_a Q^\pi(s, a) \\ \frac{\epsilon}{|A|} & \text{ov.} \end{cases}$  Pf:  $\sum \pi'(a|s) Q^\pi(s, a) \geq \sum \pi(a|s) Q^\pi(s, a)$  when  $\epsilon$  small.  
off-policy MC resample ratio  $P_\pi^\pi(s_t) \cong \frac{\pi(a_t | s_t)}{\sum_{a'} \pi(a' | s_t)}$   $\Rightarrow Q(s_t, a_t) += \frac{G - Q(s_t, a_t)}{P_\pi^\pi(s_t, a_t)}$   
( $b_k$  is soft  $b_k(a|s) > 0, \forall a, s$ )  $W^* = \frac{\pi_k(a_t | s_t)}{b_k(a_t | s_t)}$  ( $Q^\pi = E_{\tau \sim P_\pi^\pi} (G)$ )

(conv. conditions are mild and thus omitted)

TD. eval by solve Bellman eq. stochastic and online. TD(0)  $V_{(s)}^{t+1} = V_{(s)}^t + \alpha_t (r(s, a, s') + \gamma V_{(s')}^t - V_{(s)}^t)$   
 $(s, a, s') \sim \pi$ . (note is biased est. to  $V^\pi, \neq V^t$ )

MC has high var  $x$ , TD low var  $v$ . learns from complete episode and no MDP structure  $X$ . (one step) incomplete episode by bootstrapping. (exploit MDP structure)

but unbiased est.  $V_{(s)}^t$  first visit.  $n$ -step.  $(\gamma^\pi)^n V_{(s)} = E_x (\sum_{k=0}^{n-1} \gamma^k r_k + \gamma^n V_{(s_n)} | S_0 = s)$   
 $G_{(s)}^\lambda \triangleq (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{(s)}^n$   
 $G_{(s)}^n \triangleq \sum_{k=0}^{n-1} \gamma^k r_k + \gamma^n V_{(s_n)}^t$  ( $s_0 = s, a_0, \dots, s_n \sim \pi$ )  
TD( $\lambda$ ).  $= \sum_{k=1}^{\infty} (\lambda)^k (r_k + \gamma V_{(s_{k+1})}^t - V_{(s_k)}^t) + V_{(s)}^t$   
 $\lambda = 0 \rightarrow TD(0). \lambda = 1 \rightarrow 0 \rightarrow G_{(s)}^\lambda = \sum_{k=0}^{\infty} \gamma^k r_k$  MC eval i.e. TD(0). (is a bias-var tradeoff)



# 復 四

(no VFA) [or TD(co) tabular is RM to solve Beq.]  
 (VFA eval) (if  $\phi$  one-hot i.e. tabular, SGD)  
 [TD(co) is semi SGD to BE problem and SGD to BE [mean  $\pi_t$ ] in (moving target) (or  $\pi_{t+1}$ )]

地址:上海市邯郸路220号

邮编:200433

电话:65642222

and SGD to BE [mean  $\pi_t$ ] in (moving target) (or  $\pi_{t+1}$ )

(a-v version)  $Q^{t+1}(s,a) = Q^t(s,a) + \alpha_t(s,a)(r(s,a,s') + \gamma Q^t(s',a') - Q^t(s,a))$ . (RM for  $Q^{t+1} \leq \beta Q^t$ )  
 (few when model known)

SARSA. TD learn with  $\epsilon$ -greedy.  $(S_t, a_t, r_t, S_{t+1}, a_{t+1}) \sim \pi_t$   $Q$  learn  $Q^{t+1}(s,a) = Q^t(s,a) + \alpha_t(s,a)(r + \gamma \max_{a'} Q^t(s',a') - Q^t(s,a))$   
 (SARSA( $\lambda$ ) con policy) TD(co) for  $Q(S_t, a_t)$   $\pi_{t+1}(a|s_t) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} Q^t(s_t, a') \\ \epsilon & \text{otherwise} \end{cases}$  (off policy and online) (Cupd policy  $\forall \mathbb{E}_{a' \sim \pi_t} Q^t(s, a')$  is just  $\max_{a'} Q^t(s, a')$  so no resampling. ie.  $\beta Q = \beta Q$ )  
 using TD( $\lambda$ )  $\pi_{t+1}(a|s_t) = \begin{cases} 1 & \text{if } a = \arg \max_{a'} Q^t(s_t, a') \\ \epsilon & \text{otherwise} \end{cases}$  (Cupd policy  $\forall \mathbb{E}_{a' \sim \pi_t} Q^t(s, a')$  is just  $\max_{a'} Q^t(s, a')$  so no resampling. ie.  $\beta Q = \beta Q$ )

double  $Q$  learn Fact: Jensen's ineq.  $\mathbb{E} \max Q(s', a)$   
 $\rightarrow$  one choose  $a^*$  but  $\geq \max_a \mathbb{E} Q(s', a)$ . (Cupd policy  $\forall \mathbb{E}_{a' \sim \pi_t} Q^t(s, a')$  is just  $\max_{a'} Q^t(s, a')$  so no resampling. ie.  $\beta Q = \beta Q$ )  
 $(S_t, a_t, r_t, S_{t+1}) \sim b_t$  e.g.  $\epsilon$ -greedy wrt.  $\frac{Q^{t+1, A} + Q^{t+1, B}}{2}$   
 prob.  $\frac{1}{2}$   $a^* = \arg \max_a Q^{t, A}(s_{t+1}, a)$   
 $\frac{1}{2}$   $Q^{t+1, A}(s_t, a_t) = Q^{t, A}(s_t, a_t) + \alpha_t(s_t, a_t)(r_t + \gamma Q^{t, B}(s_{t+1}, a^*) - Q^{t, A}(s_t, a_t))$   
 $b^* = \arg \max_a Q^{t, B}(s_{t+1}, a)$  unbiased.  
 $Q^{t+1, B}(s_t, a_t) = \sim$

3 VFA value func approx.  $V^\pi(s) \leftarrow V(s; w)$  e.g. lin.  $\phi'(s; w)$  or  $NW$   $w(s)$ .  
 (not tabular)  $Q^\pi(s, a) \leftarrow Q(s, a; w)$ .  
 $\min_w \mathbb{E}_{s \sim D} (V(s; w) - V^\pi(s))^2$   
 SGD  $\Rightarrow w_{t+1} = w_t + \alpha_t (V^\pi(s) - V(s; w_t)) \nabla_w V(s; w_t)$

TD eval VFA perspectives  $w_{t+1} = \arg \min_w \mathbb{E}_{s \sim D} (V(s; w) - \mathbb{E}_{(s, w)} V(s; w))^2$  by MC or TD eval. (for lin. and SGD gives it.)  
 or Bellman error  $\triangleq \mathbb{E}_{s \sim D} (V(s; w) - \gamma V(s; w))^2$   $(\|V - V^\pi\|_\infty \leq \|V - \mathbb{T}^\pi V\|_\infty)$  (biased)  
 min semigrad replaces  $\mathbb{T}^\pi V(s; w)$  by  $\mathbb{T}^\pi V(s; w_t)$  also  $\mathbb{T} \cdot \mathbb{T}^\pi$  version) (MC conv. to the above min by SGD, (like IRLS) to vanish  $\nabla_w \mathbb{T}^\pi V(s; w)$ . interpreting  $Q$  learn. TD conv. to min)

learn with VFA. (a-v version) SARSA-VFA  $\nabla_{(S_t, a_t, r_t, S_{t+1}, a_{t+1})} \text{projected BE if lin.}$   
 (both follow TD target  $\beta^\pi Q$ . TDeval VFA  $\leftarrow w_{t+1} = w_t + \alpha_t (r_t + \gamma \phi(S_{t+1}, a_{t+1})' w_t - \phi(S_t, a_t)' w_t)$  with lin.  $\sim \pi_t$  (to  $\phi$ )  
 SARSA est. by samples  $\sim \pi$  and  $Q$ -learn  $\epsilon$ -greedy upd.  $\leftarrow \pi_{t+1} = \begin{cases} 1 & \text{if } a = \arg \max_{a'} \phi(S_t, a)' w_{t+1} \\ \epsilon & \text{otherwise} \end{cases}$   $\phi(S_t, a_t)$   
 use greedy  $\pi$   $\beta^\pi Q = \beta Q$ .  $Q$ -learn-VFA similar.  $\min_w \mathbb{E}_{(s, a) \sim D} (Q(s, a; w) - \beta^\pi Q(s, a; w))^2$   
 $(\pi := \beta^\pi Q)$   $\leftarrow w$  (SARSA) SGD.  
 And both upd. policy right after value para upd, which may be stab. issue in complex VFA e.g. deep  $Q$ -learn) FQI. see  $w_t =$  below. (i.e. approx VI with steps)

DQN  $\nabla$  replay buffer  $\mathcal{D}$  to capacity  $N$ ; (pin target a while (policy to upd.) and minibatch)  
 FQI with deep  $NW$  for VFA  $Q$  network  $Q(s, a; w)$ , target  $Q$  network  $Q(s, a; \tilde{w})$ ;  
 and adopts incremental learning  $\tilde{w} = w$ ; SGD iter number  $C$ ;  $k=0$ ;  $S_0$   
 by buffer) for  $t=0 \rightarrow ?$  stop if  $k=C$   
 break seq. dependence  $k+=1$   $\tilde{w} = w$   
 $(S_t, a_t, r_t, S_{t+1}) \sim b_t$  from  $S_t$ ; add into  $\mathcal{D}$ . end  $k=0$   
 (tricks: double, batch  $B$  of  $\mathcal{D}$  end  
 dueling, prioritized replay, dip..)  $w_t = \frac{1}{C} \sum_{(s, a) \in B} (r_t + \gamma \max_{a'} Q(s', a'; \tilde{w}) - Q(s, a; w)) \nabla_w Q(s, a; w)$



# 復旦大學

[06 (Tsitsiklis - Van Roy)]

地址: 上海市邯郸路220号

邮编: 200433

电话: 65642222

网址: //www.fudan.edu.cn

o PBE  $Q = \Pi \mathcal{R} Q$  where  $\text{proj. } \Pi f = \arg \min_{g \in \mathcal{Q}_\pi} \|f - g\|_{L^2(\mathcal{D})} \triangleq \arg \min_{g \in \mathcal{Q}_\pi} \mathbb{E} (f - g)^2$  (D stn. sample from  $\pi$ )  
 $V = \Pi T V$   $\mathcal{Q}_\pi \triangleq \{Q(\cdot; \omega) \mid \omega \in \mathcal{U}\}$  subsp. o6 TD(0) if lin. VFA and on policy,

above actually  $\mathbb{E}_{(s,a) \sim \mathcal{D}} (Q - \mathcal{R} Q)^2$  eval  $\pi$   $\sum \alpha_t = \infty, \sum \alpha_t^2 < \infty$ ,  
 thus  $\leq \mathbb{E}_{(s,a) \sim \mathcal{D}, s' \sim p(\cdot|s,a)} (Q - (r + \gamma \max_{a'} Q(s', a', \omega)))^2$  (i.e. TD(0)) conv. to PBE solution.

Q-learn-VFA is noisy upperbound. in tabular case (MSBE) min PBE not (PF by mt) i.e. Semi SGD's conv. condition

above means  $Q(\cdot; \omega_{t+1}) = \Pi \mathcal{R} Q(\cdot; \omega_t)$   $\Pi = I$ , solve (P)BE  $Q_{t+1} = Q_t + \alpha_t (\mathcal{R} Q_t - Q_t + e_t)$   
 but only if o6 condition ow.  $\Pi \mathcal{R}$  may not be contraction  $\nrightarrow$  PBE. (Baird convex)  $\uparrow$

model/policy-based RL. simplex para.  $\pi_s \in \Delta, \forall s, e^{f_\theta(s,a)}$  by RM  $\leq \mathbb{E} \mathbb{E}$ : all conv.  $(e_s + e_{s,a})$   
 softmax para.  $\pi_\theta(a|s) = \frac{e^{f_\theta(s,a)}}{\sum_{a'} e^{f_\theta(s,a' )}}$  (note differ  $\frac{1}{s}, (r + \gamma \max)$ )  
 $V^{\pi_\theta}(\mu) \triangleq \mathbb{E}_{s \sim \mu} V^{\pi_\theta}(s_0) = \mathbb{E}_{s \sim p^{\pi_\theta}(\cdot|s_0)} (\sum_{t=0}^{\infty} \gamma^t r_t)$ , RL  $\arg \max_{\theta} V^{\pi_\theta}(\mu)$  but VFA introduces error of  $> 0$   
 $V^{\pi}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi}} \mathbb{E}_{a \sim \pi(\cdot|s)} \mathbb{E}_{s' \sim p(\cdot|s,a)} r(s,a,s')$ , where  $d_{\mu}^{\pi}(s) \triangleq \mathbb{E}_{s_0 \sim \mu} (1-\gamma) \sum_{t=0}^{\infty} \gamma^t p(s_t = s | s_0, \pi)$  (biased but  $\xrightarrow{a.s.} Q^*$  of BE)  
 advantage func.  $A^{\pi}(s,a) \triangleq Q^{\pi}(s,a) - V^{\pi}(s)$  (discounted state visit measure)

o7 (Performance difference lem.)  $V^{\pi_1}(\mu) - V^{\pi_2}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_1}} \mathbb{E}_{a \sim \pi_1(\cdot|s)} A^{\pi_2}(s,a)$   
 Pf:  $V^{\pi_1} - V^{\pi_2} = (I - \gamma P^{\pi_1})^{-1} (T^{\pi_1} V^{\pi_2} - V^{\pi_2})$ ,  $T^{\pi_1} V^{\pi_2} = \mathbb{E}_{a \sim \pi_1(\cdot|s)} Q^{\pi_2}(s,a)$

o8 (Policy gradient)  $\nabla_{\theta} V^{\pi_\theta}(\mu) = \mathbb{E}_{s \sim p^{\pi_\theta}(\cdot|s_0)} (\sum_{t=0}^{\infty} \gamma^t Q^{\pi_\theta}(s_t, a_t) \nabla_{\theta} \log \pi_\theta(a_t | s_t))$  ( $Q^{\pi_\theta}$  here can be viewed as eval)  
 Pf:  $\nabla_{\theta} Q^{\pi_\theta}(s_0, a_0) = \gamma \mathbb{E}_{s_1 \sim p(\cdot|s_0, a_0)} \nabla_{\theta} V^{\pi_\theta}(s_1) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} (Q^{\pi_\theta}(s,a) \nabla_{\theta} \log \pi_\theta(a|s))$

PPG.  $\nabla_{\pi_s} V^{\pi}(\mu) = \frac{d_{\mu}^{\pi}(s)}{1-\gamma} Q^{\pi}(s, \cdot)$  by o8. (or  $\sim A^{\pi}(s, a)$ ) note traj. expr  $\mathbb{E}_{s \sim \pi} \sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t)$   
 upd.  $\pi^{k+1} = \arg \min_{\pi \in \Pi} \{ -\alpha_k \langle \nabla_{\pi} V^{\pi}(\mu), \pi - \pi^k \rangle + \frac{1}{2} \|\pi - \pi^k\|_2^2 \}$  if  $\sum \alpha_k = 1$ . also uses, but  $\rightarrow$  visitation meas. expr

$\Rightarrow \pi_s^{k+1} = \text{Proj}_{\Delta} (\pi_s^k + \frac{\alpha_k d_{\mu}^{\pi^k}(s)}{1-\gamma} Q^{\pi^k}(s, \cdot))$ , vs.  $O(\frac{1}{k})$  conv. ( $\frac{1}{1-\gamma}$  omit)  $(\mathbb{E}_{s,a \sim d_{\mu}^{\pi^k}} Q^{\pi^k}(s,a))$   
 Softmax PG.  $\nabla_{\theta} V^{\pi_\theta}(\mu) = \frac{d_{\mu}^{\pi_\theta}(s)}{1-\gamma} \pi_\theta(a|s) A^{\pi_\theta}(s, \cdot)$  by o8 and  $\frac{\partial \log \pi_\theta(a|s)}{\partial \theta_{s,a}} = \frac{1}{\pi_\theta(a|s)} (\mathbb{1}_{a=s} - \pi_\theta(a|s)) = \mathbb{E}_{s,a \sim \pi_\theta} Q^{\pi_\theta}(s,a) \nabla_{\theta} \log \pi_\theta(a|s)$   $\neq$  traj.  $\nabla_{\theta} (\mathbb{E}_{s,a} \sum \gamma^t r_t)$

PMA.  $\pi^{k+1} = \arg \min_{\pi \in \Pi} \{ \sum_s (-\alpha_k \langle \nabla_{\pi} V^{\pi}(\mu), \pi_s - \pi_s^k \rangle) + \frac{\alpha_k d_{\mu}^{\pi^k}(s)}{1-\gamma} D_h(\pi_s, \pi_s^k) \}$  e.g.  $h(\pi_s) = \frac{1}{2} \|\pi_s\|_2^2 \Rightarrow \pi_s^{k+1} = \text{Proj}_{\Delta} (\pi_s^k + \text{Proj}_{\Delta} \text{Ascent. } \alpha_k Q^{\pi^k}(s, \cdot))$   
 (regularizer  $h(\pi_s) = \sum_s \pi_s(a) \log \pi_s(a)$ ,  $D_h = \text{KL}$ )  
 of visi. reweighted simplex  $\arg \min_{\pi_s \in \Delta} \{ -\alpha_k \langle Q(s, \cdot), \pi_s - \pi_s^k \rangle + D_h(\pi_s, \pi_s^k) \}$  vs.  $h(\pi_s) = \sum_s \pi_s(a) \log \pi_s(a)$ ,  $D_h = \text{KL}$ .  
 Bregman divergence)  $\alpha_k \rightarrow \infty$ , PPG, PMA  $\rightarrow$  PI. Expo QAscent.  $\pi_{s,a}^{k+1} \propto \pi_{s,a}^k e^{\alpha_k Q^{\pi^k}(s,a)}$

PQA, EQA

Pf: coincides MPG.  $(\text{KL}(\pi_s^{k+1} \| \pi_s^k)) \geq 0$   
 $\rightarrow$  both  $O(\frac{1}{k})$  conv.  $= \log Z_s^k$   
 PPG, PQA conv. in finite iters.

(REINFORCE)  $\nabla$  for  $k=0 \rightarrow ?$

PG + MC eval

learn  $\pi_\theta$  via pg.

learn in  $\nabla_{\theta} Q_{\omega}$  via eval.

sample  $\mathcal{D}_k = \{s_t^i\} \sim \pi_{\theta_k}$

$G_t^i := \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}^i$

$g_k = \frac{1}{|\mathcal{D}_k|} \sum_i \sum_t \gamma^t G_t^i \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i)$

variance reduction  $Q(s, a)$  (assume  $\sum_a \pi_{\theta}(a|s) = 1$ )

actor-critic

(TD, on policy)  $(s_t, a_t, r_t, s_{t+1}, a_{t+1}) \sim \pi_{\theta}$

optimal  $\theta^* = \arg \min_{\theta} \mathbb{E}_{\pi_{\theta}} \|\nabla_{\theta} \log \pi_{\theta}\|^2$

(V func, adv ver)

cor using all the episode to upd.)

$d_t = r_t + \gamma Q(s_{t+1}, a_{t+1}; \omega) - Q(s_t, a_t; \omega)$

note  $\theta$  upd.  $\frac{d}{dt} \theta$  is est. by  $\pi_{\theta}$  sampling.

MC

$A(s_t, a_t) \approx r_t + \gamma V(s_{t+1}; \omega) - V(s_t; \omega)$

eg. MC, n-TD for  $\omega$ .

$\omega \leftarrow \omega + \alpha d_t \nabla_{\omega} Q(s_t, a_t; \omega)$

is est. by  $\pi_{\theta}$  sampling.

end

$\nabla (r_t + \gamma V(s_{t+1}; \omega) - V(s_t; \omega))$

note for  $\theta$  if not normalized (need to Proj or softmax)

$\mathbb{E}_{\pi_{\theta}} f(\theta)$

$(s_t, a_t, r_t, s_{t+1}) \sim \pi_{\theta}$

MC  $G_t$  can't as multiplier to upd.

$\max_{\theta} J(\theta)$

$d_t = r_t + \gamma V(s_{t+1}; \omega) - V(s_t; \omega)$

$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t$  (natural gradient)

$\omega \leftarrow \omega + \alpha d_t \nabla_{\omega} V(s_t; \omega)$

$\Delta \theta \propto \arg \max \{J(\theta) + \langle \nabla J(\theta), d \rangle\}$

$\theta \leftarrow \theta + \beta \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

if its argmax that's  $\leftarrow \text{KL}(\pi_{\theta} || \pi_{\theta+d}) \leq \alpha$

TRPO

comp. to  $\square 7$ , write  $V_{\theta}(\mu)$  (expand  $\pi_{\theta}(x)$  to second order) Fisher int  $F(\theta) := \mathbb{E}_{\pi_{\theta}} (\nabla_{\theta} \log \pi_{\theta}) (\nabla_{\theta} \log \pi_{\theta})'$

$\max_{\theta} V_{\theta}(\mu) := V_{\theta}^{\pi_{\theta}}(\mu) + \frac{1}{1-\gamma} \mathbb{E}_{\pi_{\theta}} A(s, a)$

in RL,  $F(\theta) = \mathbb{E}_{\pi_{\theta}} (\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t))'$

Facts:  $V_{\theta}^{\pi_{\theta}}(\mu)$  and  $V^k(\theta)$  match at  $\theta_k$

average case  $\lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi_{\theta}} \sum_{t=0}^{T-1} \dots = \mathbb{E}_{\pi_{\theta}}$

up to first derivative;

discounted case stn. dis.  $\leftarrow \mathbb{E}_{\pi_{\theta}} (\sum_{t=0}^{\infty} \gamma^t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t))'$

$V_{\theta}^{\pi_{\theta}}(\mu) \geq V_k(\theta) - \frac{2\gamma \epsilon}{(1-\gamma)^2} \max_{s,a} |A(s,a)|$

where  $w^* = \arg \min_w \mathbb{E}_{\pi_{\theta}} (\nabla_{\theta} \log \pi_{\theta}(a|s) \omega - A(s,a))^2$

$\Rightarrow \max_{\theta} V_k(\theta) \leftarrow \nabla_{\theta} V_{\theta}^{\pi_{\theta}}(\mu) (\theta - \theta_k)$

for Softmax,  $w_{s,a}^* = A^{\pi_{\theta}}(s,a) + c_s$ ,  $\pi_{\theta}^{k+1}(a|s) = \pi_{\theta}^k(a|s) \cdot e^{\frac{1}{1-\gamma} A(s,a)}$

s.t.  $\mathbb{E}_{\pi_{\theta}} \text{KL}(\pi_{\theta}^k(\cdot|s) || \pi_{\theta}(\cdot|s)) \leq \epsilon$

(pf by  $\nabla_{\theta} \log \pi_{\theta}(a|s) \omega^* = \sum_a \nabla_{\theta} \log \pi_{\theta}(a|s) (A^{\pi_{\theta}}(s,a) + c_s)$  7 stepsize.  $\sum_a \sim$

(replace by  $\mathbb{E}$  version)

$\frac{1}{2} (\theta - \theta_k)' F(\theta_k) (\theta - \theta_k) \leq \epsilon$  Same with EQA.  $\frac{\pi_{\theta}(a|s)}{\pi_{\theta}^k(a|s)} A(s,a)$

$\approx$  vPG + line search in implementation.

PPO.  $V_k(\theta) - V_{\theta}^{\pi_{\theta}}(\mu) \propto \mathbb{E}_{\pi_{\theta}} (\frac{\pi_{\theta}(a|s)}{\pi_{\theta}^k(a|s)} A(s,a))$

$\nabla$  for  $k=0 \rightarrow ?$

(surrogate of the target in small region around  $\theta_k$ , under sampling from  $\pi_{\theta}$ )

dip.  $r(\theta) \triangleq \frac{\pi_{\theta}(a|s)}{\pi_{\theta}^k(a|s)}$

$\tau_n = \{s_t, a_t, r_t, s_{t+1}\}_{t \geq 0} \sim \pi_{\theta_k}$

$1 - \epsilon \leq r(\theta) \leq 1 + \epsilon^k$

$d_t(w_k) = r_t + \gamma V(s_{t+1}; w_k) - V(s_t; w_k)$

thus  $\theta$  won't move far  $l = r(\theta_k)$   $L_{\text{clip}}(\theta, \theta_k, s, a) \triangleq$

$A^{\lambda}(s_t, a_t) = \sum_{l \geq 0} (\gamma \lambda)^l d_{t+l}(w_k)$ ;  $G^{\lambda}(s_t) = A^{\lambda}(s_t, a_t) + V(s_t; w_k)$

using PG-type method.  $\min \{r(\theta), \text{clip}(r(\theta), 1-\epsilon, 1+\epsilon) A\}$

$\theta_{k+1} = \arg \max_{\theta} \sum_{\tau_n} \sum_t L_{\text{clip}}(\theta, \theta_k, s_t, a_t)$

(lower bound conservative) actually  $G^{\lambda}(s_t, a_t)$ , e.g.  $A < 0, r > 1 + \epsilon$ , less.

$w_{k+1} = \arg \min_w \sum_{\tau_n} \sum_t (V(s_t; w) - G^{\lambda}(s_t))^2$

when it w we take average, not affecting max ...

end

in implem. use GDer. e.g. Adam,

$\theta \leftarrow \theta + \beta \sum_{i \in \mathcal{B}_t} \nabla_{\theta} L^{\text{dip}}(s_i, a_i; \theta, \theta_k)$ ;  $w \leftarrow w + \alpha (\frac{1}{|\mathcal{B}_t|}) \sum_{i \in \mathcal{B}_t} (G^{\lambda}(s_i) - V(s_i; w)) \nabla_w V(s_i; w)$



# 復旦大學

地址: 上海市邯鄲路220号

邮编: 200433

电话: 65642222

网址: //www.fudan.edu.cn

entropy regularization  $V_{\epsilon}^{\pi}(\mu) \triangleq \frac{\mathbb{E}_{s \sim d_{\mu}^{\pi}} (r(s,a,s') + \gamma V^{\pi}(s'))}{\mathbb{E}_{a \sim \pi(\cdot|s), s' \sim p(\cdot|s,a)}} = \mathbb{E}_{s \sim \mu, \pi} \left( \sum_{t=0}^{\infty} \gamma^t (r_t - \tau \log \pi(a_t|s_t)) \right)$

(like revising reward to encourage explore) Bellman op.  $T_{\epsilon}^{\pi} V(s) \triangleq \mathbb{E}_{a,s'} (r - \tau \log \pi(a|s) + \gamma V(s'))$

$\max_{\pi} T_{\epsilon}^{\pi} V \Rightarrow \pi(a|s) \propto e^{\frac{Q_{\epsilon}^{\pi}(s,a)}{\epsilon}}$ .  $Q_{\epsilon}^{\pi}$  def the same,  $V_{\epsilon}^{\pi}(s) = \mathbb{E}_{a \sim \pi(\cdot|s)} (Q_{\epsilon}^{\pi}(s,a) - \tau \log \pi(a|s))$

$\square 1'$   $V_{\epsilon}^{\pi_1}(\mu) - V_{\epsilon}^{\pi_2}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_1}(s)} (\mathbb{E}_{a \sim \pi_1(\cdot|s)} A_{\epsilon}^{\pi_1}(s,a) - \tau \text{KL}(\pi_1(\cdot|s) || \pi_2(\cdot|s))) \mathbb{E}_{a \sim \pi_2(\cdot|s)} A_{\epsilon}^{\pi_2}(s,a) = 0$ ,  $A_{\epsilon}^{\pi} \triangleq Q_{\epsilon}^{\pi} - V_{\epsilon}^{\pi} - \tau \log \pi$

$\square 2'$   $T_{\epsilon} = \max_{\pi} T_{\epsilon}^{\pi}$ .  $\exists \pi^* \text{ s.t. } V_{\epsilon}^{\pi^*} = V_{\epsilon}^*$ .  $T_{\epsilon}^{\pi_1} V_{\epsilon}^{\pi_2}(s) - V_{\epsilon}^{\pi_2}(s)$  since  $\pi^* \propto e^{\frac{Q^*}{\epsilon}}$ ,  $A_{\epsilon}^* = 0$ , v.s.a.

Soft+PI  $\pi_{k+1}(\cdot|s) = e^{\frac{Q_{\epsilon}^{\pi_k}(s,\cdot)}{\epsilon}}$ .  $\tau \rightarrow 0$ ,  $\pi_{\epsilon}^* \rightarrow \pi_0^*$ .  $V(s) = \mathbb{E}_{s' \sim p(\cdot|s,a)} (r(s,a,s') + \gamma V(s')) - \tau \log \pi(a|s)$ , v.s.a.

$\square 3'$   $\nabla_{\theta} V_{\epsilon}^{\pi_0}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_0}} (\mathbb{E}_{a \sim \pi_0(\cdot|s)} \nabla_{\theta} \log \pi_0(a|s)) V_{\epsilon}^* \rightarrow \max_a Q^* = V^*$ .  $V = V_{\epsilon}^*$ ,  $\pi = \pi_{\epsilon}^*$ . Pf: Softmax Lem.

(OSS) also for Softmax para. the same. e.g. ent. sm. PG  $\pi_{k+1} \propto \pi_k e^{\frac{\gamma d_{\mu}^{\pi_k}(s)}{1-\gamma} \pi_k(a|s) A_{\epsilon}^{\pi_k}(s,a)}$

DPG.  $a = \pi_{\theta}(s)$ .  $\square 3''$   $\nabla_{\theta} V_{\epsilon}^{\pi_0}(\mu) = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{\mu}^{\pi_0}} (\nabla_{\theta} \pi_{\theta}(s) \cdot \nabla_a Q^{\pi_0}(s, \pi_{\theta}(s))) \sim \frac{\gamma}{1-\gamma} A_{\epsilon}^{\pi_0}(s,a)$  i.e.  $\propto \pi_k$  s.s.a.

Pf:  $\nabla_{\theta} (Q^{\pi_0}(s_0, \pi_{\theta}(s_0))) = \nabla_{\theta} \pi_{\theta}(s_0) \nabla_a Q^{\pi_0}(s_0, a) + \gamma \int \nabla_{\theta} V_{\epsilon}^{\pi_0}(s_1) p(s_1|s_0, a = \pi_{\theta}(s_0)) ds_1$

DDPG. (a-c  $Q(s,a|w) \rightarrow Q^{\pi_0}(s,a)$ , by FQI. b need add noise to  $\pi_0$ , replay buffer for dependence.)

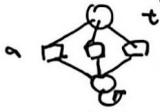
5 online planning. (from now on, no  $\pi(\cdot|\cdot)$ )  
 receding horizon planning (depth  $d$ ), MCTS.

Dyna-Q.  
 (combines model-based methods on estimated model i.e. simulation experience to improve efficiency of usage.)  
 data



building incrementally by e-e tradeoff.

MAB.  $\nabla$  SARSA/Q-learn  
 for  $a_t \sim \epsilon_t$ -greedy  $Q^t$  (or  $b$ )



$Q(a) = Q^t(a_t) + \alpha_t (r_t - Q^t(a_t))$   
 (Hoeffding ineq.)

regret  $R_T \triangleq \sum_{t=1}^T (M_t^* - M_{a_t})$ ,  $M_a \triangleq \mathbb{E}_{r \sim D_a}(r)$   
 for evaluating algo. (note for offline (orade) regret can be 0.)

$\{a_1, \dots, a_k\}$   $X_k$  iid. sub-Gauss with  $v$ .  $\mathbb{P}(|\frac{1}{n} \sum_{i=1}^n (X_i - \mu)| \geq t) \leq 2e^{-\frac{nt^2}{2v^2}}$   
 ( $\mathbb{E} e^{\lambda(X-\mathbb{E}X)} \leq e^{\frac{\lambda^2 v^2}{2}}, \forall \lambda \in \mathbb{R}$ )

algorithms:  $\langle \text{UCB} \rangle$  esp. Gauss. bd.  
 explore-first.

e-greedy.  $\bar{\mu}_t(a) = \frac{\sum_{i=1}^t r_i \mathbb{1}_{\{a_i=a\}}}{\sum_{i=1}^t \mathbb{1}_{\{a_i=a\}}}$  (same as SARSA/Q-learn)  
 $\begin{cases} 1-\epsilon + \frac{\epsilon}{k} & \text{when } d_t = \text{sth.} \\ \frac{\epsilon}{k} & \dots \end{cases}$

( $N$  rounds each  $a$ )  $\bar{\mu}_a = \frac{1}{N} \sum_{t=1}^N r_{a,t}$   
 then phase II  $T - nk$  rounds  $\hat{a} = \arg \max_a \bar{\mu}_a$ .

$\mathbb{E} R_T \leq T^{\frac{2}{3}} O(k \log T)^{\frac{1}{3}}$  if  $\nu = (\frac{T}{k})^{\frac{2}{3}} O(\log T)^{\frac{1}{3}}$

$\epsilon_t = (\frac{k}{t})^{\frac{1}{3}} \log t^{\frac{1}{3}}$

UCB.  $\arg \max_a \bar{\mu}_t(a) + \text{radius}_t(a)$

by o/b,  
 $|\bar{\mu}_t(a) - \mu(a)| \leq v \cdot \sqrt{\frac{2}{n(a)} \log \frac{2}{\delta}}$  with Prob.  $1-\delta$

$v^2$ : variance proxy  
 e.g.  $\sigma^2$  Gauss.  $\text{supp} \subset [-M, M]$   $\mathbb{P}(X \geq a) \leq \frac{\text{inf}_{t \geq 0} \text{supp} \mathbb{P}(X \geq t) e^{-ta}}{M^2}$   
 Hoeffding ineq. can be derived Chernoff bound.

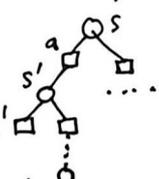
typically,  $\delta = \frac{1}{T}$ ,  $v_{\text{UCB}} = \bar{\mu}_t(a) + C \sqrt{\frac{\log nt}{2n_t(a)}}$   
 as radius

$\mathbb{E} R_T \leq O(\sqrt{kT \log T})$  if  $\delta = \frac{2}{T}$ ,  $\{ \mu_a \}_{a=1}^k$  obey a priori dis.,  
 $\forall t \leq T$ .  
 Bayes Bandits. select  $a$  by  $\mathcal{P}(\arg \max_{a'} \mu_{a'} = a | \mathcal{H}_{t-1})$

MCTS (multi-step and conducts UCB search)

$n(s, a) + 1$  i.e.  $Q(s, a) + C \sqrt{\frac{\log n(s)}{n(s, a)}}$   
 $Q(s, a) = \frac{1}{n(s, a)} (r(s, a, s') + \gamma V(s')) - Q(s, a)$

if  $a$  is chosen at  $s$  ( $V(s')$  is assumed)  
 note after this  $s'$  is not leaf.  
 (until  $s'$  meets termination conditions)



select-expand-propagate (traceback  $\gamma$ -ahead (use  $V(s')$  or simulate) e.g.  $r_0 + \gamma r_1 + \dots + \gamma^n V(s_n)$ )

Thompson Sampling.  $\mathcal{H}_{t-1} = \{a_1, r_1, a_2, r_2, \dots, a_{t-1}, r_{t-1}\}$   
 (sample  $(\mu_a)_{a=1}^k \sim \mathcal{P}(\cdot | \mathcal{H}_{t-1})$  choose the even ind. setting i.e. largest.  
 $\bar{\mu}_t(a) \sim \mathcal{P}(\mu_a | \mathcal{H}_{t-1}), a=1, \dots, k$   
 $\Rightarrow$  optin.  $O(\sqrt{kT \log T})$

AlphaGo Lee: supervised or RL get a learned net for policy fast rollout (net as simulator) BP. upd  $Q(s, a), \pi(s, a)$ .  
 Selection  $Q(s, a) + C \frac{P(s, a)}{1+n(s, a)}$   
 Expand  $\geq n_0$   
 Evaluation  $V(s) = \lambda V_0(s) + (1-\lambda) \otimes$

Zero:  $Q + C \cdot \frac{P(r(s))}{1+n(s, a)}$  (PUCT)  
 same net has  $p, v$  two heads  
 learn from MCTS (use  $\pi$  from MCTS search guided by  $p$ , use  $r(s, a)$  as  $\pi$  till end)  
 $L = (\delta - v)^2 - \pi \log p + c \|o\|^2$

# 6 RL for LLMs. token-level MDP.

RM for RLHF  $\left\{ \begin{array}{l} \text{rule-based} \\ \text{discriminative} \\ \text{generative} \end{array} \right.$  (deterministic) e.g. safety or structural constraints, e.g. math\_verify coding  
 compare or rank based on human preference, e.g. RM<sub>o</sub>(x, y) ∈ ℝ

DPO. from max<sub>θ</sub> V<sub>θ</sub>(D) by LLM e.g. CoT.  
 $\Rightarrow \pi_{\theta}^* = \frac{e^{\beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)}}}{Z_{\beta}(x)}$

under BT model,  $P(y_c > y_r | x) = \sigma(\beta \log \frac{\pi_{\theta}(y_c|x)}{\pi_{\text{ref}}(y_c|x)} - \beta \log \frac{\pi_{\theta}(y_r|x)}{\pi_{\text{ref}}(y_r|x)})$   
 Bradley-Terry model  $\log \frac{P(y_c > y_r)}{1 - P(y_c > y_r)} = r_{\theta}(x, y_c) - r_{\theta}(x, y_r)$   
 classification head decoders [x, y]

thus, min DPO obj.  $\min_{\theta} \mathbb{E}_{(x, y_c, y_r) \sim \mathcal{D}} -\log \sigma(\beta \log \frac{\pi_{\theta}(y_c|x)}{\pi_{\text{ref}}(y_c|x)} - \beta \log \frac{\pi_{\theta}(y_r|x)}{\pi_{\text{ref}}(y_r|x)})$  ie. min  $\mathbb{E} -\log \sigma(r_{\theta}(x, y_c) - r_{\theta}(x, y_r))$   
 max<sub>θ</sub> V<sub>θ</sub>(D) =  $\mathbb{E}_{x \sim \mathcal{D}} (\mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} r_{\theta}(x, y) - \beta \text{KL}(\pi_{\theta}(\cdot|x) || \pi_{\text{ref}}(\cdot|x)))$

(off-line friendly, need high-qa. data) PG.  $\nabla_{\theta} V_{\theta}(\mathcal{D}) = \mathbb{E}_{x, y} \sum_{t=0}^{T-1} \nabla \log \pi_{\theta}(y_t | x, y_{<t})$   
 GRPO. sample  $\{y_i\}_{i=1}^k$  for x, PPO/TRPO  $\mathbb{E}_{x, y \sim \pi_{\theta}(\cdot|x)} \sum_{t=0}^{T-1} \frac{\pi_{\theta}(y_t | x, y_{<t})}{\pi_{\text{ref}}(y_t | x, y_{<t})} A_{\theta}^k(y_t, x, y_{<t})$   
 est. by  $A_{t,i}^{\pi_{\theta}} = \frac{r_{\theta}(x, y_i) - \text{mean}(r_{\theta}(x, y_{1:k}))}{\text{std}(r_{\theta}(x, y_{1:k}))}$

note  $A_{t,i}^{\pi_{\theta}} = Q_{\theta}^{\pi_{\theta}}(x, y_{<t}, y_t) - V_{\theta}^{\pi_{\theta}}(x, y_{<t})$  but vanilla PPO use GAE(λ)  $\xrightarrow{\text{est.}} A_{t,i}^{\pi_{\theta}}$  and intermediate rewards is difficult for even RM trained.  
 $\approx Q_{\theta}^{\pi_{\theta}}(x, y_{<t}, y_t) - V_{\theta}^{\pi_{\theta}}(x, y_{<t})$  sufficiently close to  $Q_k$ . e.g. dipped obj.  $\min_{\theta} \sum_{t=0}^{T-1} \text{clip}(\frac{r_{\theta}(x, y_t)}{\pi_{\text{ref}}(y_t|x)}, 1 \pm \epsilon) A_{t,i}^{\pi_{\theta}}$   
 unbiased  $\frac{r_{\theta}(x, y) - \text{mean}(r_{\theta}(x, y_{1:k}))}{\text{std}(r_{\theta}(x, y_{1:k}))}$

(note  $V_{\theta}^{\pi_{\theta}}$  here needn't include  $\beta \text{KL}$  term)

at step k, sample  $\{x_i\}_{i=1}^M \sim \mathcal{D}, \{y_{ij}\}_{j=1}^M \sim \pi_{\theta}(\cdot|x_i)$   
 total policy loss  $\mathcal{L}_k(\theta) = \frac{1}{NM} \sum_i \sum_j \mathcal{L}_k(\theta | x_i, y_{ij})$   
 $\mathcal{L}_k(\theta) = -\mathbb{E}_{x, y \sim \pi_{\theta}(\cdot|x)} \sum_{t=0}^{T-1} \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\text{ref}}(a_t | s_t)} (r(x, y) - \beta \sum_{t'=0}^{T-1} \log \frac{\pi_{\theta}(a_{t'} | s_{t'})}{\pi_{\text{ref}}(a_{t'} | s_{t'})})$   
 has  $\nabla_{\theta} \mathcal{L}_k(\theta_k) = -\nabla_{\theta} J(\theta_k)$

REINFORCE. (by MC)  $\mathcal{L}_k(\theta | x_i, y_{ij}) = -\sum_{t=0}^{T-1} \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\text{ref}}(a_t | s_t)} r_{\theta}(x_i, y_{ij})$ , unbiased but high variance.  
 reduce var: mod.  $-\beta \sum_{t=0}^{T-1} \dots \rightarrow -\beta \sum_{t=0}^{T-1} \dots$  by noting that  $\nabla_{\theta} J(\theta) = \mathbb{E}_{x, y} \sum_{t=0}^{T-1} \nabla \log \pi_{\theta}(a_t | s_t)$   
 (but  $r(x_i, y_{ij})$  also has var.)

A-C. (like in GRPO)  $V_{\theta}(s_t) \triangleq \mathbb{E}_{y \sim \pi_{\theta}(\cdot|x)} r(s_t, y), Q_{\theta}(s_t, a) = V_{\theta}(s_t, a)$   
 then A-C trains a model to approx  $Q^{\pi_{\theta}}$  (or  $V^{\pi_{\theta}}$ ) in  $\mathcal{L}_k$ .  
 $(r(x_i, y_{ij}) | s_t^{(ij)}, a_t^{(ij)}) \rightarrow \hat{Q}(s_t^{(ij)}, a_t^{(ij)})$  no var. but has bias in  $Q$ 's approx.  
 i.e.  $\mathbb{E} \nabla_{\theta} \mathcal{L}_k(\theta_k | x_i, y_{ij}) \neq -\nabla_{\theta} J(\theta_k)$

$\hat{\mathcal{L}}_k(\theta) = \frac{1}{NM} \sum_i \sum_j \left( \sum_{t=0}^{T-1} \frac{\pi_{\theta}(a_t^{(ij)} | s_t^{(ij)})}{\pi_{\text{ref}}(a_t^{(ij)} | s_t^{(ij)})} \hat{Q}(s_t^{(ij)}, a_t^{(ij)}) \right)$   
 Lists of estimation of  $\hat{Q}_k(s_t, a_t)$ .  
 (note if  $\min_{\theta} \mathcal{L}_k(\theta)$  we may enlarge  $\pi_{\theta}(a_t | s_t)$  for all sampled tokens  $a_t$  (ow. softmax or  $\sum=1$  would adjust)  
 since  $\hat{Q}_k^{\pi_{\theta}}(s_t, a_t) > 0$  and  $\{a_t^{(ij)} | s_t^{(ij)}\}$  not covers all vocabulary  
 all" reduce a baseline  $\hat{\mathcal{L}}_k(\theta | x_i, y_{ij}) = -\sum_{t=0}^{T-1} \left( \frac{\pi_{\theta}(a_t^{(ij)} | s_t^{(ij)})}{\pi_{\text{ref}}(a_t^{(ij)} | s_t^{(ij)})} \cdot (\hat{Q}(s_t^{(ij)}, a_t^{(ij)}) - b(s_t^{(ij)})) \right)$   
 and choose  $b(s_t) = V^{\pi_{\theta}}(s_t)$   
 for  $\mathbb{E}_{a_t \sim \pi_{\theta}(\cdot | s_t)} Q^{\pi_{\theta}}(s_t, a_t) = V^{\pi_{\theta}}(s_t)$

7 Meta RL.  $\min_{\theta, \psi} \mathbb{E}_{T \sim P(T)} \mathcal{L}(D_T^{\text{test}}, \theta')$  s.t.  $\theta' = f_{\psi}(D_T^{\text{tr}}, \theta)$ .  $f_{\psi}$   $\begin{cases} \text{gradient-based} \\ \text{context-based} \\ \text{(recurrence-)} \\ \text{e.g.} \end{cases}$

Grad. MAML  $f_{\psi}(D, \theta) = \theta - \alpha \nabla_{\theta} \mathcal{L}(D, \theta)$  (or multistep)  
 $\psi = \{\alpha_i\}$ ; or  $\phi$ .

$\nabla$  sample  $T_i \sim P(T)$   
 for  $i$  sample  $\{\tau_{ij}^{\theta}\}_{j=1}^k$  from  $f(D_{T_i}, \theta)$  in  $T_i$ .  
 $\theta_i' = \theta - \alpha \nabla_{\theta} \mathcal{L}_i(\theta)$  ( $\mathcal{L}_i(\theta) := \frac{1}{k} \sum_{j=1}^k \mathcal{L}_i(\tau_{ij}^{\theta})$ )  
 end sample  $\{\tau_{ij}^{\theta_i'}\}_{j=1}^m$  from  $f(D_{T_i}, \theta_i')$  in  $T_i$ .  
 upd.  $\theta := \beta \nabla_{\theta} \sum_i \mathcal{L}_i(\theta_i')$ . (FO:  $\frac{\partial \theta_i'}{\partial \theta} = I - \alpha \frac{\partial \mathcal{L}_i(\theta)}{\partial \theta}$ )

MAESN.  $\theta, \{(\mu_i, \sigma_i)\}_i$ .  $\begin{matrix} \uparrow \\ \square \\ \uparrow \\ s \end{matrix} \begin{matrix} \uparrow \\ \square \\ \uparrow \\ s \end{matrix} \begin{matrix} \uparrow \\ \square \\ \uparrow \\ s \end{matrix} \begin{matrix} \uparrow \\ \square \\ \uparrow \\ s \end{matrix} \sim \mathcal{N}_i$   
 $\max_{\theta, \mu_i, \sigma_i} \sum_{T_i} \mathbb{E} (\sum_{t \sim \pi(\cdot|s_t, \theta)} r_t(s_t, a_t) - \sum_{t \sim \pi(\cdot|s_t, \theta)} \text{KL}(\mathcal{N}(\mu_i, \sigma_i) \parallel \pi(\cdot|s_t, \theta)))$   
 s.t.  $\mu_i' \sim \mathcal{N}(\mu_i, \sigma_i)$   
 $\mu_i' = \mu_i + \alpha \nabla_{\mu_i} \mathbb{E} (\sum_{t \sim \pi(\cdot|s_t, \theta)} r_t(s_t, a_t))$   
 $\sigma_i' = \sigma_i$

Context. POMDP belief bcs.

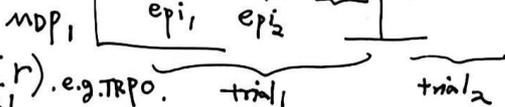
$b_{t+1}(s') \propto \Omega(o|s', a) \sum_s P(s'|s, a) b_t(s)$ . eg. using TRPO.

(action  $a$  and observe  $o$ ) but still train

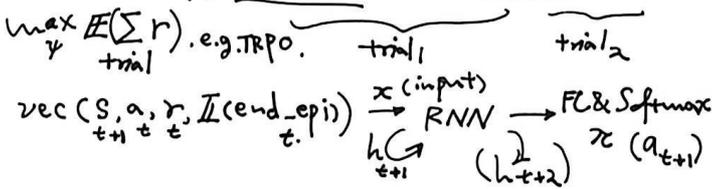
(On-policy mostly, PEARL. but Pearl can be off-policy)

RL<sup>2</sup>

$\psi$  in RNN.  $h_0 \rightarrow h_{T_1} \rightarrow h_{T_1+T_2} \rightarrow \dots$  at episode-level (though intro recur. structure in deep RL)  
 eg. GRU.  $s_0 \rightarrow s_{T_1} \rightarrow s_0 \rightarrow s_{T_2} \rightarrow \dots$



better when trial contains less episodes.



7-4  $(\max_{\theta} \mathbb{E}_{x \sim \mathcal{D}} Q(x, \pi_{\theta}(x)))$  use max (greedy implicitly).  $\hat{\theta} \leftarrow \text{Polyak}$

DDPG (DPG + DQN) (soft upd i.e.  $\tau\theta + (1-\tau)\hat{\theta}$ ) TD3 (clipped double Q)  $r + \gamma \min\{Q^A(s', a'), Q^B(s', a')\}$   
 (delayed  $\pi$  upd.) actor and target net upd./C step.

sample  $\pi_{\theta} + \mathcal{N}_{\epsilon} = a_t$  into  $\mathcal{D}$ .  $\{sars'\}_t$   
 thus  $\theta$  upd. just like  $\mathbb{E} \max_{a'} Q$  which leads to higher est. of  $Q_w$ . (target i.e. previous  $w$ ) to avoid Q over-fitting.  $a_t = \pi_{\theta}(s_t) + \text{clip}(\mathcal{N}(0, \sigma^2))$

soft Q learn.  $Q^{t+1}(s, a) = Q^t(s, a) + \alpha (s, a) (r + \gamma \tau \log \sum_{a'} e^{\frac{Q^t(s', a')}{\tau}} - Q^t(s, a))$ . temp.  $\tau$  adjusting: target ent.  $\mathcal{L}_{\tau} \triangleq \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi_{\theta}(s)} -\tau (\log \pi_{\theta}(a|s) + \mathcal{H}_{\theta})$ .

SAC (off-policy i.e. replay  $\mathcal{D}$ )  $\alpha$  target  $\tilde{Q} = r + \gamma (\min_{A, B} Q_w^A(s', V_{\pi}^B(s')) - \tau \log \pi_{\theta}(a|s'))$ . where  $\pi_{\theta}(a|s) = \frac{e^{-\beta Q_w(s, a)}}{\sum_{a'} e^{-\beta Q_w(s, a)}}$ .  $a = \mu_{\theta}(s) + \sigma_{\theta}(s) \cdot \epsilon$ ,  $\epsilon \sim \mathcal{N}(0, 1)$ .

VariBAD

(online)

$$z_\phi(z|z_t) = \mathcal{N}(\mu_t, \sigma_t^2)$$

e.g. GRU.

$$P_\theta(\tau|z) = \prod_{t=1}^T P_\theta(s_{t+1}|s_t, a_t, z)$$

VAE 练习 task repae.

$$\pi_\psi(a_t|s_t, z_\phi(z|z_t))$$

$$\mathcal{L} = \mathbb{E}_{(\omega, \psi, \gamma)} \left[ \mathcal{J}(\psi, \phi) + \lambda \sum_{t \in \mathcal{E}} \text{ELBO}_t(\phi, \theta) \{x, z\} \right]$$

PEARL

(offline)

where RL  $\mathcal{J} = \mathbb{E}_{\tau \sim \pi_\psi, z_\phi} \left[ \sum_{t \in \mathcal{E}} \gamma^t r_t + \alpha H(\tau, \psi) \right]$

$$z \sim z_\phi(z|\text{context})$$

(from  $\pi_\theta \mathcal{T}_i$ )

$$\pi_\theta(a|s, z)$$

$\theta, \omega$  by SAC

$$\psi \text{ by } \nabla_{\phi} \sum_i (\mathcal{L}_i^{\text{critic}} + \beta \text{KL}(z_\phi \| p(z)))$$

VAE:

$$x \xrightarrow{\text{en.}} z \xrightarrow{\text{de.}} \hat{x}$$

$$z_\phi(z|x) \quad p_\theta(x|z)$$

(AE deterministic  $z, \hat{x}$ )  
recon.  $\mathcal{L} = \|x - \hat{x}\|^2$

(if not have this term it comes to AE)

$$\mathcal{L} = -\mathbb{E}_{z \sim z_\phi(x)} (\log p_\theta(x|z)) + \beta \cdot \text{KL}(z_\phi(z|x) \| p(z))$$

(recon.: neg. log. likelihood)  $\beta > 1$  info. bottleneck  
eg. Softmax  $\rightarrow$  Cross Ent. (ie. disentangle)

(grad thru  $z_t$ ) Gaussian  $\rightarrow$  MSE.  $\beta < 1$  relieve posterior collapse

$$z_\phi(z|x) = \mu_\phi(x) + \sigma_\phi(x) \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

$$\mathcal{L}_{\text{recon.}} = \mathbb{E}_{\epsilon} (\log p_\theta(x | \mu_\phi(x) + \sigma_\phi(x) \epsilon))$$

Gen. model.  $(\max_x) \log p_\theta(x)$

$$\text{note eq. here} \quad = \log \int p_\theta(x|z) p(z) z_\phi(z|x) dz$$

$$\log p_\theta(x) \geq \mathbb{E}_{z_\phi(z|x)} \log p_\theta(x|z) - \text{KL}(z_\phi|p(z))$$

$$\geq \text{ELBO} \quad (\text{max ELBO})$$

(ie.  $p_\theta(z|x)$  is difficult to calculate due to  $\frac{1}{z}$  then use  $z_\phi$  to approx) ie.  $\min \text{KL}(z_\phi(z|x) \| p_\theta(x, z))$

intuition:  $\nabla_{\theta} \log p_\theta(x) = \frac{\int p(z) \nabla_{\theta} \log p_\theta(x|z) dz}{p_\theta(x)}$  but  $\int p_\theta(z|x) \nabla_{\theta} \log p_\theta(x|z) dz$

if no encoder we still can not get  $p_\theta(z|x)$ , note intro.  $z_\phi(z|x)$  we remain two KL (sum of  $\text{KL}(z_\phi|p_\theta(z|x)) - \text{KL}(z_\phi|p(z))$ )

ie {Randomness} (separa. explicitly)

'search in recurrence'  $\Delta z$  每步  $\pi_\theta$  时 fix.

评价: MAESV  $\pi_\theta(a|s, z)$   $z \sim \mathcal{N}(\mu_i, \sigma_i^2)$

没有 history coherent. (grad in 总)

是所有的 hist. 作用.

RL<sup>2</sup>

episode 不 引导 for cost. (upd. hist./hidden)

GRU RNN 改 LLM + Mload.

无 randomness (不对  $\mathcal{T}_i$  作分布 inference)

task infer. i.e.  $\pi^*$  与 task infer. disentangle (C.18)

$z_\phi(z|s)$

PEARL

encoder  $z_\phi(z|\text{Context})$  for infer.

排列对称.  $c_i = \{s_i, a_i, r_i, s'_i\}$ .

VariBAD

VAE construct  $\tau$ .

$\pi_\psi(a_t|s_t, \mu_{\phi,t}, \sigma_{\phi,t})$  本质又回到了 recurrence.

(只是 hidden 的提取 use recons. 进行)

(Me to Q/learn also GRU  $\pi_\theta(a_t|s_t, z_t)$ )

explicit: 1/2. low score 选 1/2 task (based on TD3 offline.) match

$z$  use GRU but also explicit  $\pi_\theta$  upd.

'critic net'