Value Function Methods - PG methods suffer from high variance when estimating Po J(4) - AC methods introduce fitted value iteration to reduce the variance in the All estimate of Do J(b) Why do we need the actor To? . can we get rid of it? - anatomy of RL algorithme with value thurching approx. Bt Qp fit a model to estimate return generate sænsles O(nu joolicy) improve policy e.g. use a  $(\varepsilon)$ -greedy policy:  $T \in T = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ , it  $\alpha = \operatorname{argmax} \mathcal{R}_{0}^{T}$ 

-/-

- cannot bit 
$$V_{p}^{T}$$
, time we can't importe  
the policy on line  
- argmax  $A^{T}(s, a) = \arg(a \times Q^{T}(s, a))$   
=> fit the Q tunction  $\rightarrow Q_{p}^{T}$   
- Algorithm: Filled Q-iteration  
a. collect data set (repeat S time)  
 $2(si, ai, ri, si)]_{i=1}^{i=1}$   
using some policy  
(Q-learning is off-policy)  
2. compute target:  
 $(x \in Y) = r(si, ai) + Y \max(Q_{p}^{T}(s', a)) - \frac{Y}{1})^{2}$   
s. do supervised regression  
 $Y \in \arg(a) = \frac{Z}{1} ||Q_{p}^{T}(s', a)| - \frac{Y}{1}||^{2}$   
assumed  
 $work \cdot inp$   
- subtordinery (G)

- Q-learning is off-policy  
- Q-learning is off-policy  
trajectoriel  
but every transition s ->s'l's indep.  
token into account  
issue: it doeen't work!  
- Algonitum: loulice Delearning (withing Q-learning)  
1. take action as using some policy  
I observe (strag. rj. str): single transition  
2. til = v(strag. rj. str): single transition  
3. 
$$\varphi \in \operatorname{argmin} \frac{1}{2} || Q_{\varphi}(strag) - ytill2$$
  
no run 2 !  
uote: argmin  $\frac{1}{2} || Q_{\varphi}(strag) - ytill2$   
 $= R_{\varphi} \frac{1}{2} || Q_{\varphi}(strag) - ytill2$   
 $= (R_{\varphi} Q_{\varphi}(strag)) + (Q_{\varphi}(strag) - ytill2) = 0$   
is treaked !  
to some 1





- 5-



-> new transitions sampled using, e.g. E-greedy policy w.r.t. Qp Note: we don't immediately train on the most recent deta point (maybe we never use it for training) . transitions sampled <u>ind</u> from butter -6-

→ O-iteration w/ replay butter:  
1. collect data set ? (sj:aj:rj:sj:)}.  
nsing some policy; add it to butter B  
2. sample a batch from B uniternly  
3. 
$$\varphi = \varphi - \Delta \sum R_p Q_p(sj:aj) \times j = 3di$$
  
 $\times \left[ Q_p(sj:aj) - [r(sj:aj) + y max Q_p(sj:aj)] \right]$   
 $\varphi = \frac{1}{2} + \frac{Q - network}{2} + \frac{1}{2} + \frac{Q - network}{2} + \frac{1}{2} + \frac{1$ 

-7-

-8-

1. save tanget Qui-net params y'exp  
2. collect dataset ? (sjiqi, M; sji)? usig  
( some policy; add it to butter B;  
( N 33. sample a batch from B uniformly  
( N update Qui-network param  
y e p-d D; 
$$\pm \sum 1/Q_p (sjiqi) -$$
  
- (r(sjiqi) + f max B/s', g')]/?  
· inner loop of steps 3 Dy regresses on  
stationary tangets di  
- algorithm: DRN:  
(Q-learning with a replay butter l' tanget Quet)  
1. take aj using some policy, e.g. E-gierdy with  
observe (sjiqi, n', sj'); odd it to butter B.  
2. sample minibadels ? (sjiqi, rj, sj'); from  
B uniformly  
3. compute: y = r(sjiqi) + y max Q/sj, q')  
using tanget Qy-network  
y is indep. of y  
-g.

-10-

avoid evaluating a & function at an  
action a' which the arguax of the same  
Q function  
onse by & target net By:  
standard Q-learning: 
$$y = r + p R_{p'}(s', arguar R_{p'}(s', d))$$
  
duble Q-learning with gates the  
overestimation of Q-values!  
- Practical implemention: "best prochice"  
- untwork architecture for  $R_{p'}(s, a)$   
i)  $s \longrightarrow DNN_{p} \longrightarrow Q(s, a)$   
i)  $s \longrightarrow DNN_{p} \longrightarrow Q(s, a)$   
i)  $s \longrightarrow DNN_{p} \longrightarrow Q(s, a)$   
i)  $R(s, a_{1AI}) = a_{1AI}$   
- Q-learning takes time to stabilize:  
· test your code on a simple take first  
e.g. cartpole with a simple state space  
mountain car, etc.

-14-

-15-