

Key challenge: explore vs. exploit

① "Optimism under uncertainty"

$N(s, a)$  # times visited  $(s, a)$

Learning

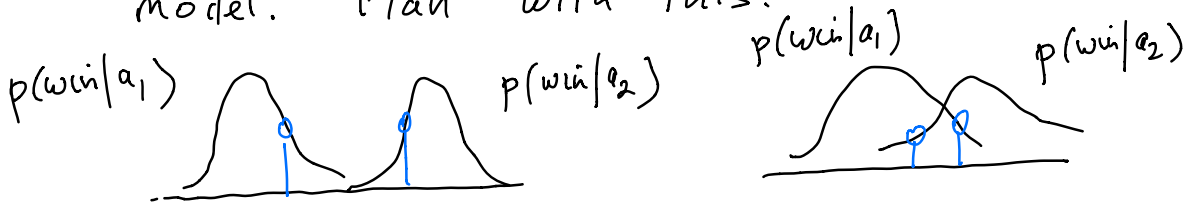
if  $N(s, a)$  is small, assume next state has high reward.

→  $R_{max}$

Plan with the "optimistic" model. Either do well, or learn new things.

② "Posterior sampling" (Thompson sampling)

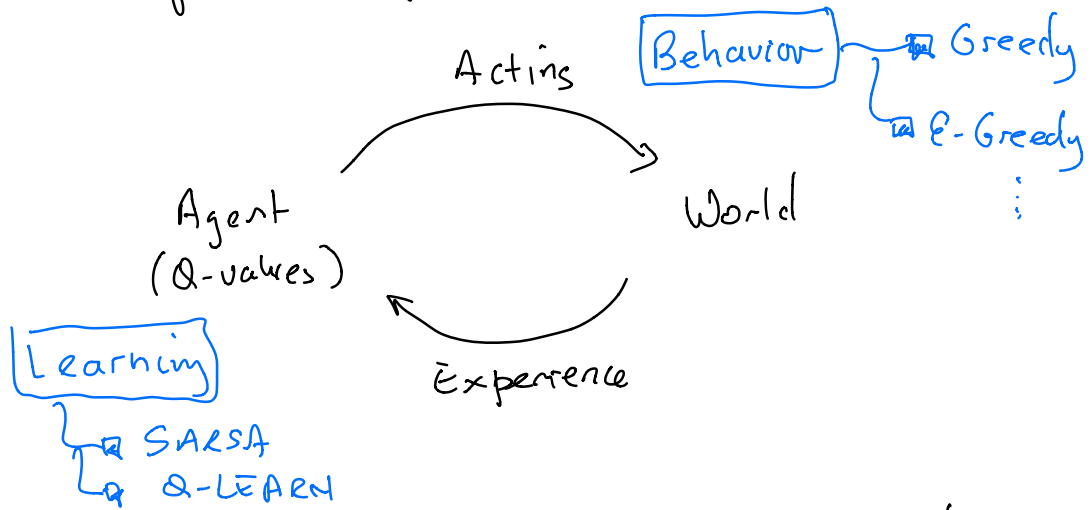
Maintain a posterior (Dirichlet) on transition  $p(\cdot | s, a)$  & Sample to get a model. Plan with this.



play  $a_2$

sometimes play  $a_1$

## B Model free Reinforcement Learning (Action-Value)



Recall

$$Q^*(s, a) = r(s, a) + \gamma \sum_{s'} p(s'|s, a) \max_{a'} Q^*(s', a')$$

$$\pi^*(s) \in \arg \max_a Q^*(s, a)$$

SARSA

$$Q(s, a) \leftarrow Q(s, a) + \alpha_{\tau}(s, a) [r + \gamma Q(s', a') - Q(s, a)]$$

On-policy. Learns Q-values that correspond with behavior.

Q-LEARN

$$Q(s, a) \leftarrow Q(s, a) + \alpha_{\tau}(s, a) [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

TD error

Off-policy. Learns Q-values that correspond to the Bellman equations ( $Q^*$ ).

Require  $\alpha_t(s,a) \leq 0$  unless  $(s,a)$  visited at  $t$ .

Theorem Q-LEARN converges to  $Q^*$  as  $t \rightarrow \infty$   
as long as

$$\textcircled{1} \sum_t \alpha_t(s,a) = \infty, \text{ all } s,a$$

$$\textcircled{2} \sum_t \alpha_t^2(s,a) < \infty, \text{ all } s,a$$

For SARSA to converge to  $Q^*$ , also need

$\textcircled{3}$  Behavior is "greedy in the limit"  
(eventually follows Q-values)

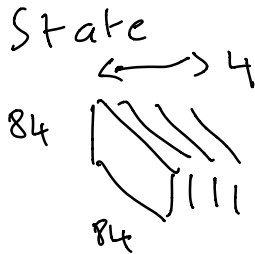
### Notes

1. For  $\sum_t \alpha_t = \infty$ , need to visit each  
 $(s,a)$  infinitely often [role of  $\epsilon$ -greedy]

2. For  $\sum_t \alpha_t^2 < \infty$ , need to reduce  
learning rate; e.g.,  $\alpha_t(s,a) = \frac{1}{N_t(s,a)}$

3. For "Greedy in limit",  
typical  $\epsilon_t(s) = \frac{1}{N_t(s)}$   $\square$

# [C] Deep-Q Networks · state 18 Q-values



Raw input, gray scaled, pixel map.  $84 \times 84$  pixels. }  $84 \times 84 \times 4$   
 4 most recent frames } 2

## Idea one

Tabular representation fails!

Parameterize  $Q(s, a; \underline{w})$ , use differentiable deep network (CNN)

## Idea two

Gradient descent on TD-error

$$\underline{w} \leftarrow \underline{w} - \frac{1}{2} \alpha_t \nabla_{\underline{w}} \left[ r + \gamma \max_{a'} Q(s', a'; \underline{w}_{old}) - Q(s, a; \underline{w}) \right]^2$$

$$+ \alpha_t \left( r + \gamma \max_{a'} Q(s', a'; \underline{w}_{old}) - Q(s, a; \underline{w}) \right) \nabla_{\underline{w}} Q(s, a; \underline{w})$$

↑  
"frozen network"

## Idea three

Experience replay

Put  $(s, a, r, s')$  into replay buffer, and do minibatch gradient descent steps.

(Re-use experience)

## ① Model-free: Policy Learning

Adopt a differentiable policy  $\pi_{\theta}(a|s)$ , parameters  $\theta$ . Learn directly!

✓ Policy might be simpler than Q-function

✓ Handle continuous actions ✓ Introduce expert knowledge

How to train?

$$\theta \leftarrow \theta + \alpha_t \nabla_{\theta} J(\theta)$$

Data histories  $h = (s, a, r, s', a', \dots)$  performance ( $\approx V^{\pi}$ )

$$\text{Reward } r(h) = \sum_t r_t \quad \text{and} \quad J(\theta) = \mathbb{E}_h [r(h)]$$

$$\text{Likelihood } \mu_{\theta}(h) = \prod_t \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_h [r(h)] = \nabla_{\theta} \int_h \mu_{\theta}(h) r(h) dh = \int_h r(h) \nabla_{\theta} \mu_{\theta}(h) dh$$

Challenge  $\nabla_{\theta} \mu_{\theta}(h)$  seems to depend on  $p(\cdot | s, a)$

Useful identity  $\nabla_{\theta} \mu_{\theta} = \mu_{\theta} \frac{1}{\mu_{\theta}} \nabla_{\theta} \mu_{\theta} = \mu_{\theta} \nabla_{\theta} \ln \mu_{\theta}$

$$= \int_h r(h) \mu_{\theta}(h) \nabla_{\theta} \left[ \sum_t \ln \pi_{\theta}(a_t | s_t) + \underbrace{\sum_t \ln p(s_{t+1} | s_t, a_t)}_{\text{doesn't depend on } \theta} \right] dh$$

$$= \mathbb{E}_h \left[ r(h) \sum_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t) \right]$$

SGD update. Given history  $h = (s, a, r, s', a', \dots)$

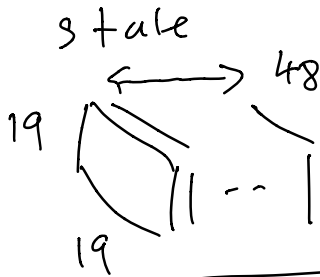
$$\theta \leftarrow \theta + \alpha r(h) \sum_t \nabla_{\theta} \ln \pi_{\theta}(a_t | s_t)$$

Can also use mini-batches

Notes <sup>(1)</sup> this is an "on policy" method because it's updating the policy that it's using to behave in the world.

<sup>(2)</sup> can be combined with "Actor-Critic" methods to speed-up + stabilize learning.

# (E) Alpha Go

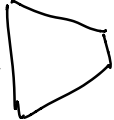


48 hand crafted, per-position features.

## Supervised learn

$\pi_{SL}$   $\pi_{fast}$

$\pi_{SL}$  Train deep net to predict expert moves

state   $19 \times 19 + 1$   
softmax (pass)

13 conv. layers

log likelihood loss

3 weeks, 50 GPUs, 30M expert pos.

$\pi_{fast}$  Train linear softmax model w/ log likelihood loss + hand coded features

## Reinforcement L

$\pi_{\theta}$   $V_{\theta}$

$\pi_{\theta}$  state   $19 \times 19 + 1$

Initialize net. to  $\pi_{SL}$

Reward +1 win  
-1 loss 0 draw

1M games vs self play vs. earlier versions

Policy gradients.

1 day, 50 GPUs, 1M games

$V_{\theta}$  state   $[0, 1]$   
 $P(win)$

Supervised learning -

Target: win or loss

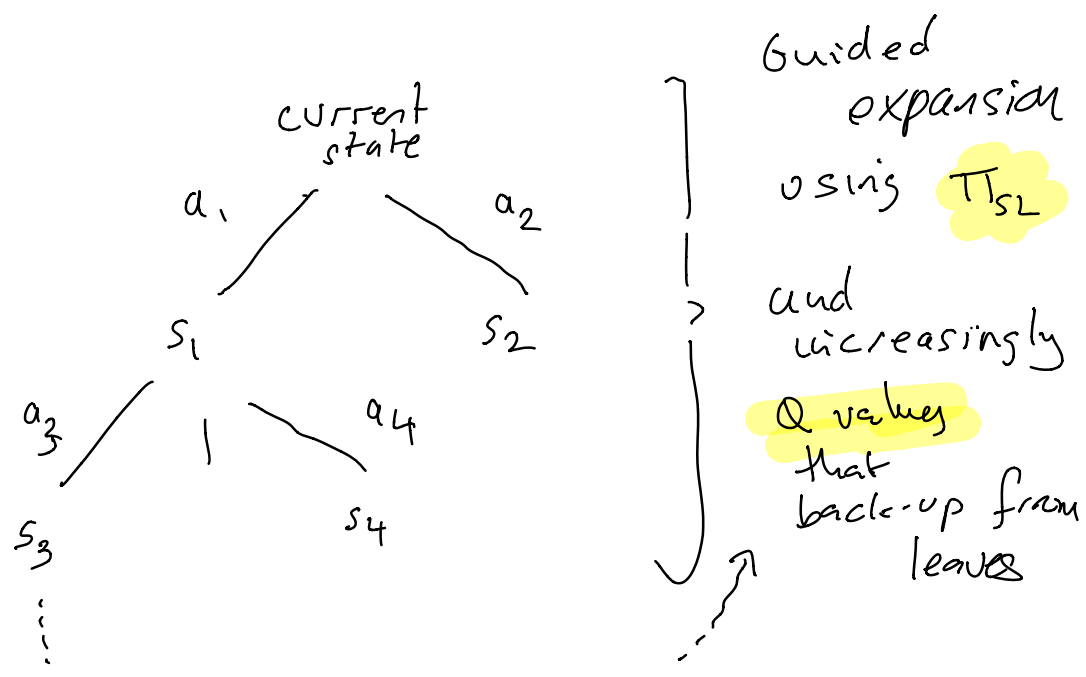
coming from Monte Carlo

eval. of  $\pi_{\theta}$ . 1 week 50 GPUs  
30M positions

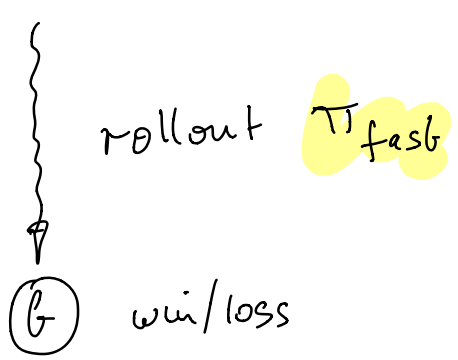


# How to play?

Guided Monte Carlo Tree Search  
 ( 48 CPUs, 8 GPUs, 40 threads !!! )



$(s_L)$  leaf  $V(s_L) = \frac{1}{2} V_{\theta}(s_L) + \frac{1}{2} G$



Didn't use  $\pi_{\theta}$ .  
 Found  $\pi_{SL}$  for guiding expansion was more useful.