# LECTURE 20

## Alpha Go

state $= \left( \begin{array}{ccc} \boxed{\cdot} & \boxed{\quad} & \boxed{\text{all } \%} \\ \bullet \ 19 \times 19 & \circ \ 19 \times 19 & \text{whose turn?} \end{array} \right)$

$S$

Action $\in \boxed{\cdot} \ + $ pass
$a \quad 19 \times 19$ positions

policy: $P(a|s)$



prob

$+ $ pass

$19 \times 19$

value: $v(s) = E_{policy}(z|s)$



$\boxed{s} \quad \begin{array}{c} z \\ +1 \\ -1 \\ +1 \\ 0 \end{array}$ average

$Q(s,a) = E_{policy}(z|s,a)$



$\boxed{s} \ (a) \quad \begin{array}{c} +1 \\ +1 \\ -9 \\ -9 \end{array}$ average

## Deep RL



$a$
prob.
softmax $\quad 19 \times 19 + 1$

$v$

conv-ResNet

$s \quad 19 \times 19 \times 3$

<u>supervised</u>:      human teacher

$$\Delta \sigma \propto \frac{\partial \log P_\sigma(a|s)}{\partial \sigma}$$

<u>RL</u>:

         RL                    supervised (fixed)

$s$ —•$\xrightarrow{P_\rho}$ (a) → $s'$ —○$\xrightarrow{P_\sigma}$ (a') → $s''$ —•$\xrightarrow{P_\rho}$ (a'') → .... → $z$

            self $P_\rho$

$$\Delta \rho \propto \frac{\partial}{\partial \rho} \log P_\rho(a|s) \cdot z \begin{matrix} +1 \\ -1 \\ 0 \end{matrix}$$
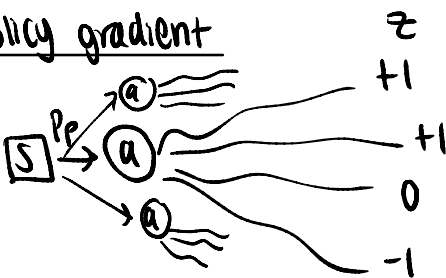
                    ↓
                  reinforcement

stochastic gradient ascent on $E(z)$
    ↳ try to max expected payout

<u>policy gradient</u>



                   $z$
                   +1
                   +1
                   0
                   -1

- if action good, ↑ prob. of action
- if action bad, ↓ prob. of action

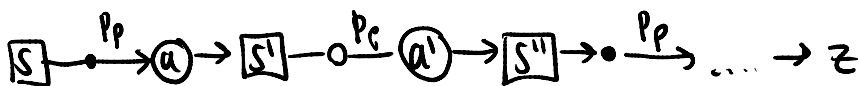$$\max E_{P_\rho(a|s)}[z] = \sum_a z \cdot P_\rho(a|s)$$

<span style="color:blue">average over payoffs of all trajectories</span>

$$\frac{\partial}{\partial \rho} E(z) = \sum_a z \frac{\partial}{\partial \rho} P_\rho(a|s)$$

$$= \sum_a \left[ z \cdot \frac{\partial}{\partial \rho} \log P_\rho(a|s) \right] \times P_\rho(a|s)$$

$$= E_{P_\rho(a|s)} \left[ z \cdot \frac{\partial}{\partial \rho} \log P_\rho(a|s) \right]$$

                       ↾
                       $\Delta \rho$

<u>learn $V_\theta(s)$</u>

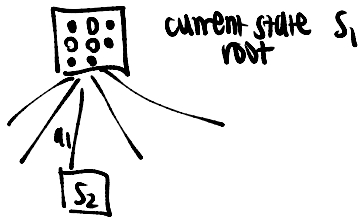$s$ —•$\xrightarrow{P_\rho}$ (a) → $s'$ —○$\xrightarrow{P_\sigma}$ (a') → $s''$ →•$\xrightarrow{P_\rho}$ .... → $z$

$$\Delta \theta \propto -\frac{\partial}{\partial \theta}(z - V_\theta(s))^2$$
         ↳ gradient descent

play the game

current state $S_1$
root

$a_1 \sim P_\rho(a|S_1)$ (impulse/reflex)

OR

$a_1$ to max $V_\theta(S_2)$

(desire)

fast thinking

$S_2$

---

$S_1$ current state

$a_1$   $a_1$   $V_\theta(S_2)$

$S_2$

$a_2$

$a_1$

$S_3$

slow thinking in imagination think ahead

$S_{20}$

$V_\theta(S_{20})$

reduce depth

Which one is more accurate?

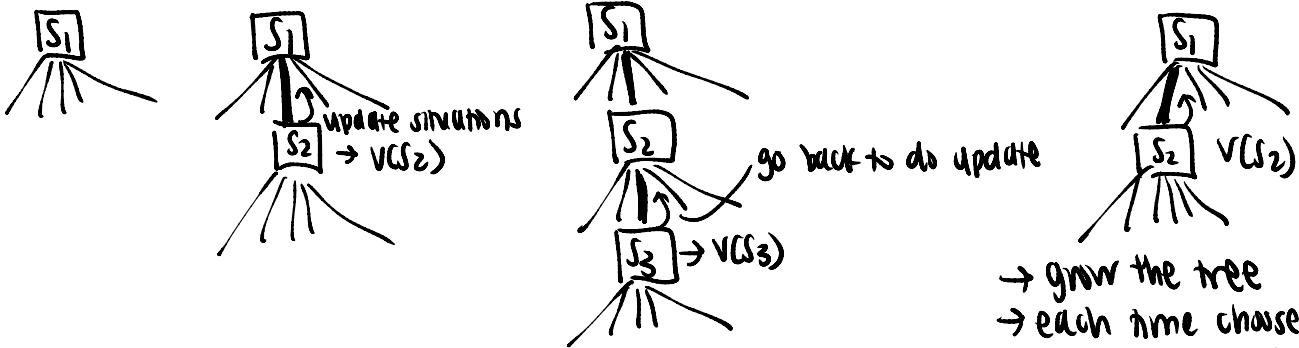$V_\theta(S_{20})$ because

closer to the end of game
⇒ more accurate
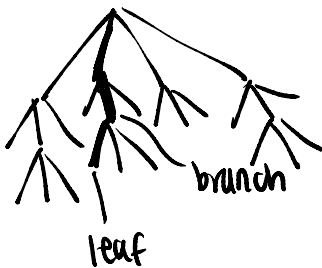    estimate of average score

monte-carlo tree search

choose best $a_1$
by max $Q(S_1, a)$
or sample $a \propto N(S,a)$

---

MCTS

$S_1$

$S_1$
$S_2$ → $V(S_2)$
update situations

$S_1$
$S_2$
$S_3$ → $V(S_3)$
go back to do update

$S_1$
$S_2$  $V(S_2)$

→ grow the tree
→ each time choose
   $a_1$ based on criteria
→ can't see full game
   tree, randomly select
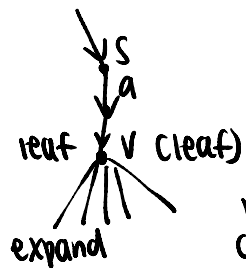   some branches

---

step 1: selection: go down a branch

branch

leaf

**Step 2:** expand

leaf

**Step 3:** back up

go back the branch

## back-up

$S$
$a$

leaf $V$ (leaf)

expand
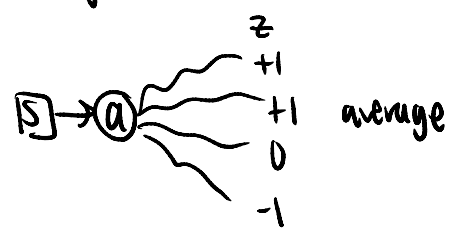
more accurate (close to end)

$N(S,a) = N(S,a) + 1$  (# of visits)
$W(S,a) = W(S,a) + V$  (total score)
$Q(S,a) = W(S,a) / N(S,a)$

(average over leaf

$S \rightarrow a$

$z$
$+1$
$+1$  average
$0$
$-1$

## Selection

$S$
$a$

choose $a$ by max
$Q(S,a) + U(S,a) \rightarrow$ exploration

exploitation

↓
uncertainty

$U(S,a) \propto \dfrac{P_\sigma(a|s)}{N(S,a)+1}$

policy according to supervised learning

gives randomness

reduce breadth

— balance btwn what know best and exploring

① train 2 policy networks — one supervised / one RL
② train value network
③ MCTS

## Alpha Go Zero
- No human data (CP6)
- Self play by MCTS

- subconscious vs. conscious
  ↓
  MCTS

( more explicit thinking



(P.V.)          (P.V.)
  ↓               ↓
$\boxed{S}$ —MCTS→ $(a)$ → $\boxed{S'}$ —MCTS→ $(a')$ → ... → $z$

$P(a|S)$        $P(a'|S')$      $V(S')$

$V(S)$

## RL in general
Markov decision process (MDP)

$\boxed{S_0}$ → ... → $\boxed{S_t}$ —$a_t$→ $\boxed{S_{t+1}}$ —$a_{t+1}$→ $\boxed{S_{t+2}}$ → ...

$r_t$ + $r_{t+1}$ + ...

(no $r_t$ in Go
$r_t$ in Atari)

## World Model
① Dynamics model:
  Deterministic $S_{t+1} = F(S_t, a_t)$

  Stochastic $P(S_{t+1} | S_t, a_t)$

② Reward model:
  Deterministic: $r_t = r(S_t, a_t)$
  Stochastic: $P(r_t | S_t, a_t, S_{t+1})$

  Reward to go: $G_t = r_t + r_{t+1} + ...$
  Goal: max $E(G_t)$          **max cumulative reward**

    policy: $P(a|S)$
    value: $V(S), Q(S,a)$

## Model-based
  known (1),(2), plan by eg MCTS          **Go**
              play out in imagination

## Model-free
  Do not know (1),(2)
  play out in real environment,          **swimming, riding bicycle**
  practice in real world