Policy : function $\pi$: observed info $\rightarrow$ action
position

$v(\pi)$ = expected reward following policy $\pi$   prob of win for QFL

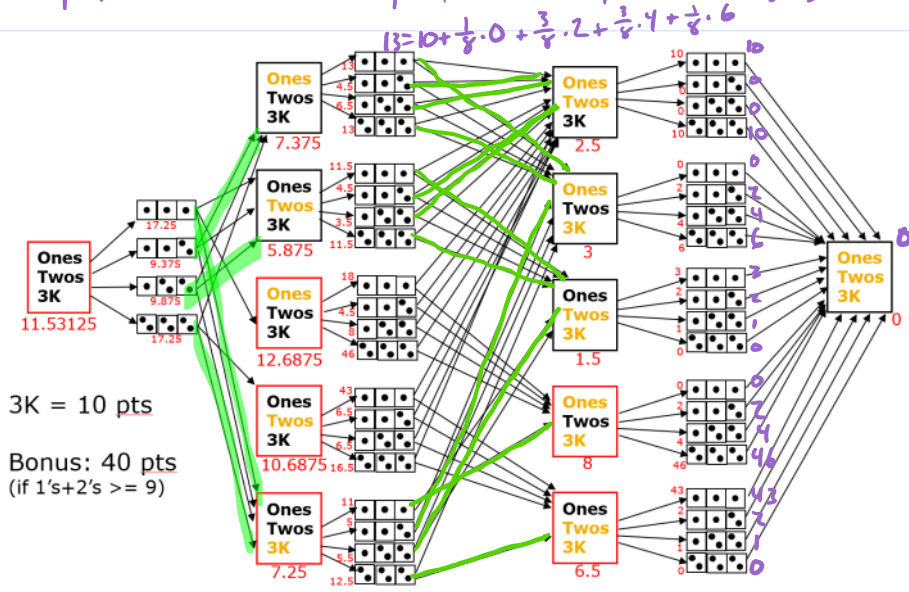sample : noisy but quick

exact : possibly slow

$\downarrow$

$$V_\pi(s) = \begin{cases} \text{value determined by game} & \text{if } s \text{ is terminal} \\ \displaystyle\sum_{s'} \underbrace{P(s \xrightarrow{\pi(s)} s')}_{\substack{\text{prob selected} \\ \text{action results} \\ \text{in state } s'}} \cdot \Big( \underbrace{R(s \xrightarrow{\pi(s)} s')}_{\substack{\text{reward earned} \\ \text{by selected action}}} + V_\pi(s') \Big) \end{cases}$$

future reward earned
from new state

for a finite game, can compute exactly!

red = values under optimal policy
purple = values under policy defined by choices highlighted in green

$13 = 10 + \frac{1}{8}\cdot 0 + \frac{3}{8}\cdot 2 + \frac{3}{8}\cdot 4 + \frac{1}{8}\cdot 6$



3K = 10 pts

Bonus: 40 pts
(if 1's+2's >= 9)

value of taking action $\hat{a}$ in state $s$     $= Q(s,a)$

$$V_a(s) = \sum_{P_a(s\to s')\neq 0} P_a(s \to s')\big(R_a(s \to s') + V(s')\big)$$

special case of MDP
$(\gamma = 1 \text{ finite})$

$$\pi_{OPT}(s) = \underset{a}{\text{argmax}} \; V_a(s) \qquad \text{optimal action is one that maximizes value}$$

$$V_{OPT}(s) = V_{\pi_{OPT}(s)}(s) \qquad \text{value of state = value of optimal action from that state}$$

## Learning Values

need discount to
avoid getting trapped
in cycles
— reward 1

observe $s \xrightarrow{a} s' \quad {=\pi(s)}$ giving $\underset{\wedge}{\text{immediate}}$ reward $R(s, s', a)$

estimate of
future reward from $s'$

get discounted future reward $\gamma \cdot V^\pi(s')$ following policy $\pi$

$\rightarrow$ discount
(maybe 1 for short
finite game)

sample of $V^\pi(s) = R(s, s', a) + \gamma V^\pi(s')$

update estimate $V^\pi(s) \leftarrow V^\pi(s) + \underbrace{\alpha}_{\text{learning rate}} \Big( \underbrace{R(s,s',a) + \gamma V^\pi(s') - V^\pi(s)}_{\text{error}} \Big)$

$$= (1-\alpha) V^*(s) + \alpha \Big( R(s,s',a) + \gamma V^\pi(s') \Big)$$

$$\max_a \sum_{s \to s'} P(s \xrightarrow{a} s') \cdot \Big( R(s \xrightarrow{a} s') + V(s') \Big)$$

$\epsilon$-greedy: w/ prob $\epsilon$ choose random
$\quad\quad\quad 1-\epsilon$ choose action $a$

but requires knowledge of model:
$P(s \xrightarrow{a} s')$ and $R(s \xrightarrow{a} s')$

$$Q(s,a) = \text{estimate of expected discounted future reward when choosing action } a \text{ in state } s$$

$$V(s) \approx \max_a Q(s,a)$$

initialize $Q(s,a) = \begin{cases} \widehat{R(s)} & \text{for terminal } s \\ 0 & \text{otherwise} \end{cases}$

*if you know this*

while not done

Sample: $r + \gamma \underline{V(s')}$

$\max_a Q(s',a)$

$s \leftarrow s_0$ *initial state*

episode

while $s$ not terminal

choose action $a$ ←

$\varepsilon$-greedy is good
prob $\varepsilon$ choose $a$ randomly uniformly
$1-\varepsilon$ choose $\text{argmax}_a Q(s,a)$

observe transition $(s,a,r,s')$

*immediate reward* $(R(s,s',a))$

update $Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a'} Q(s',a') - Q(s,a))$

*learning rate (can ↓ as episodes ↑)*

*surprise (error)*

— learn a fxn that approximates $Q(s,a)$

## Linear Approximator

Define features of states and possibly actions
↳ yards-to-1st-down / plays left

can ignore action for features of state only

$f_1(s,a)$      on pace to earn upper bonus

$f_2(s,a)$      is chance unused

⋮      both LS, SS unused

     upper category a' is unused

$$Q(s,a) = w_1 \cdot f_1(s,a) + \cdots + w_n f_n(s,a)$$

learn the weights

In state $s$

Choose action $a$ using exploit/explore policy

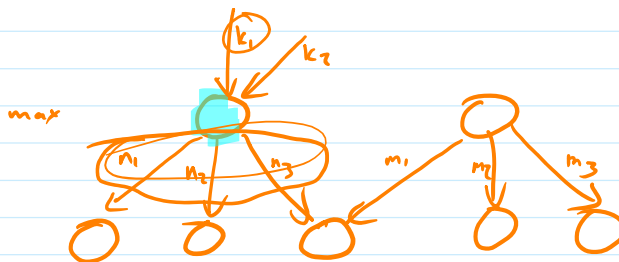Observe transition $(s, a, r, s')$

Update    ~~$Q(s,a) \leftarrow Q(s,a) + \alpha (r + \max_{a'} Q(s',a') - Q(s,a))$~~

$$w_i \leftarrow w_i + \alpha \left( r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right) \cdot f_i(s,a)$$



$k_4$

$(p, c)$

$$r + \sqrt{\frac{2 \ln n_1 + n_2 + n_3}{n_i}}$$