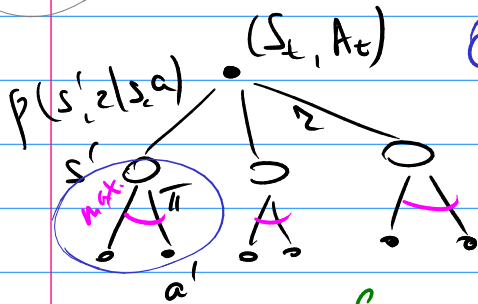
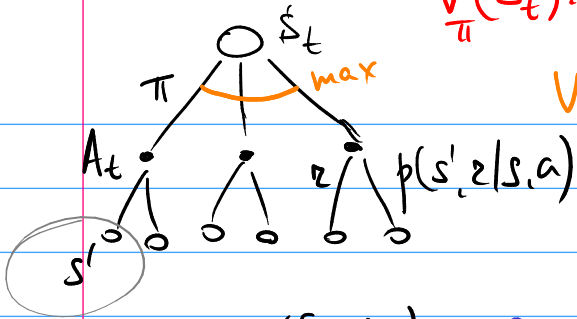


V

$$V_{\pi}(S_t) := \sum_a \pi(a|S_t) \sum_{s',z} p(s',z|s,a) (z + \gamma V_{\pi}(s'))$$

$$V_*(S_t) = \max_a \sum_{s',z} p(s',z|s,a) (z + \gamma V_*(s'))$$

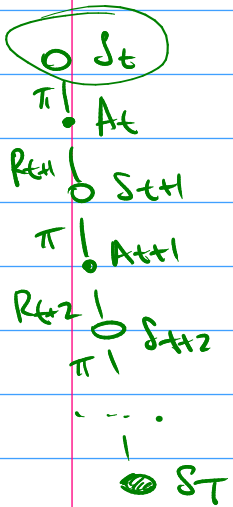
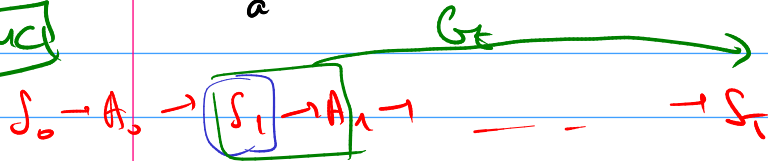


$$Q_{\pi}(S_t, A_t) = \sum_{s',z} p(s',z|s,a) (z + \sum_{a'} \pi(a'|s') Q_{\pi}(s', a'))$$

$$Q_*(S_t, A_t) = \sum_{s',z} p(s',z|s,a) (z + \max_{a'} Q_*(s', a'))$$

$\uparrow$   
 $p(s',z|s,a)$

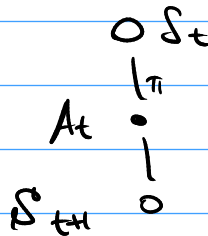
[MC]



[TD]

TD(0)

$$V_{\pi}(S_t) := V_{\pi}(S_t) + \alpha (R_{t+1} + \gamma V_{\pi}(S_{t+1}) - V_{\pi}(S_t))$$



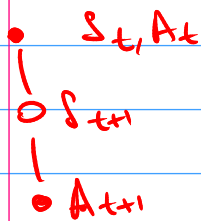
On-policy TD

Sarsa

off-policy TD

Q-learning

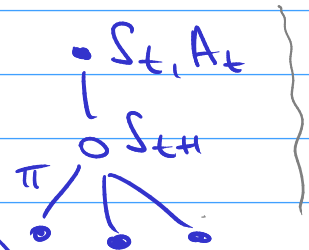
$$Q(S_t, A_t) := Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$



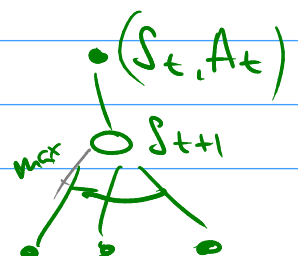
Expected Sarsa

$$Q := Q +$$

$$\alpha (R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a) - Q)$$



$$Q := Q + \alpha (R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q)$$



TD:  $\hat{G}_t = R_{t+1} + \gamma V(S_{t+1})$  n-step

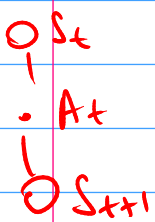
$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} = R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} =$$

$$= R_{t+1} + \gamma R_{t+2} + \gamma \left( \sum_{k=0}^{\infty} \gamma^k R_{t+k+3} \right) = \dots$$

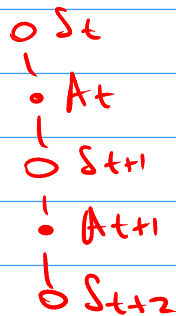
$$G_{t:t+2} = R_{t+1} + \gamma R_{t+2} + \gamma^2 V(S_{t+2})$$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

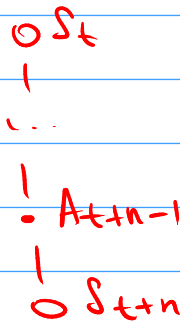
TD(0)



2-step TD



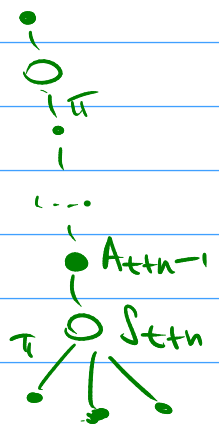
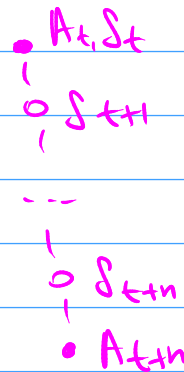
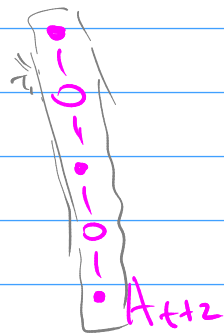
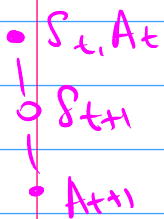
n-step TD



$$V_{\pi}(S_t) = V_{\pi}(S_t) + \alpha (G_{t:t+n} - V_{\pi}(S_t))$$

(on-policy)  
expected Sarsa

n-step Sarsa:



$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha (G_{t:t+n} - Q(S_t, A_t))$$

$\pi$  - on-policy

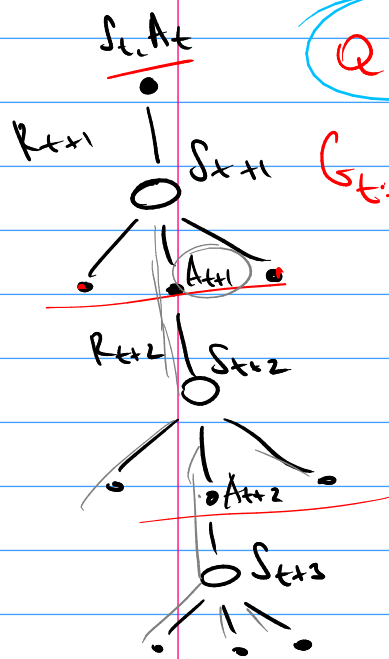
$\pi'$  - no-begreene arena

$$V(S_t) = V(S_t) + \alpha \cdot \underbrace{p_{t:t+n-1}}_{\text{probability}} (G_{t:t+n} - V(S_t))$$

$$p_{t:t+h} = \prod_{k=t}^{t+h} \frac{\pi(A_k | S_k)}{\pi'(A_k | S_k)}$$

$$Q(S_t, A_t) \approx Q(S_t, A_t) + \alpha \underbrace{p_{t+1:t+n}}_{\text{memory}} [G_{t:t+n} - Q(S_t, A_t)]$$

n-step off policy  
sarsa



$$Q(S_{t+1}, A_{t+1}) = Q(S_t, A_t) + \alpha (G - Q(S_t, A_t))$$

Tree backup

$$G_{t:t+3} = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q(S_{t+1}, a) +$$

$$+ \gamma \pi(A_{t+1}|S_{t+1}) (R_{t+2} + \gamma \sum_{a \neq A_{t+2}} \pi(a|S_{t+2}) Q(S_{t+2}, a))$$

$$+ \gamma \pi(A_{t+2}|S_{t+2}) (R_{t+3} + \gamma \sum_a \pi(a|S_{t+3}) Q(S_{t+3}, a))$$

Tree

$$G_{t:t+n} = R_{t+1} + \gamma \sum_{a \neq A_{t+1}} \pi(a|S_{t+1}) Q(S_{t+1}, a) + \gamma \pi(A_{t+1}|S_{t+1}) \cdot G_{t+1:t+n} =$$

$$= R_{t+1} + \gamma \underbrace{\sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a)}_{\tilde{V}(S_{t+1})} + \gamma \pi(A_{t+1}|S_{t+1}) (G_{t+1:t+n} - Q(S_{t+1}, A_{t+1}))$$

$$= R_{t+1} + \gamma \tilde{V}(S_{t+1}) + \gamma \pi(A_{t+1}|S_{t+1}) (G_{t+1:t+n} - Q)$$

OFF

$$G_{t:t+n} = R_{t+1} + \gamma p_{t+1} (G_{t+1:t+n} - Q) + \gamma \tilde{V}(S_{t+1})$$

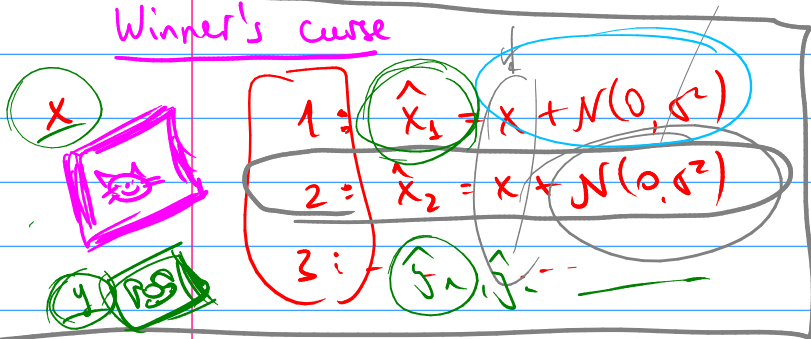
$Q(\sigma)$

$\sigma_{t+1}$

$$G^\sigma = R_{t+1} + \gamma (\sigma_{t+1} p_{t+1} + (1 - \sigma_{t+1}) \pi(A_{t+1}|S_{t+1})) (G_{t+1:t+n} - Q) + \gamma \tilde{V}(S_{t+1})$$

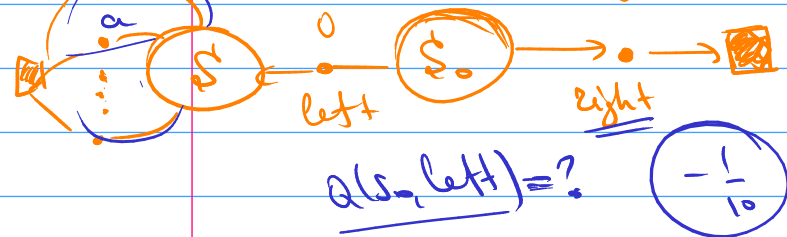
# Maximization bias, double learning

## Winner's curse



$$N(-\frac{1}{10}, 1)$$

$$Q(s_0, \text{right}) = 0$$



$$Q(s_0, \text{left}) = ? \quad -\frac{1}{10}$$

$$\max_a Q(s, a) > 0$$

## Double Q-learning:

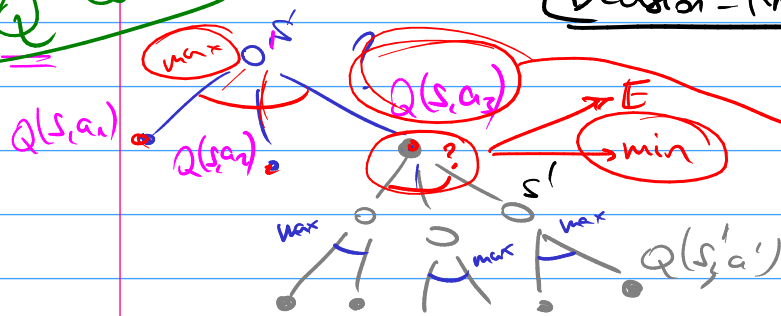
- space monkey

$$Q_1(s_t, a_t) := Q_1(s_t, a_t) + \alpha (R_{t+1} + \gamma \max_a Q_2(s_{t+1}, a) - Q_1(s_t, a_t))$$

$$Q = Q^*$$

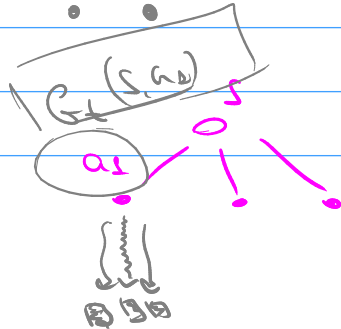
## Decision-time planning

### ① Heuristic search



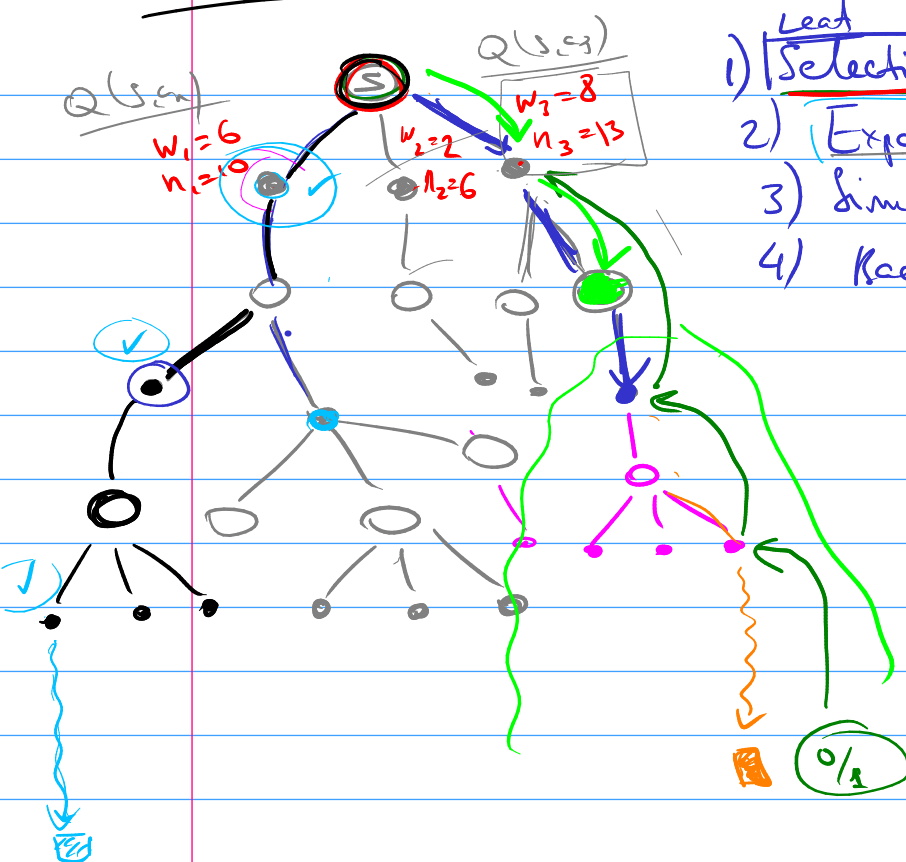
$$\alpha Q(s, a) + (1 - \alpha) \hat{Q}(s, a)$$

### ② Rollouts

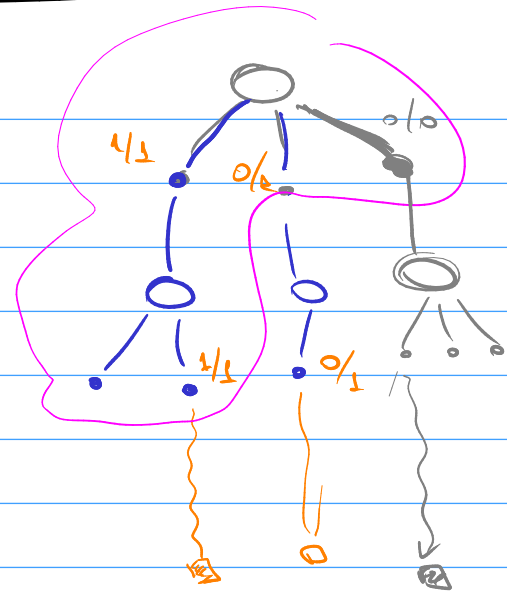


Tesauro - Hoedon

③ MCTS - Monte Carlo Tree Search



- 1) <sup>Leaf</sup> Selection
- 2) Expansion
- 3) Simulation
- 4) Recovery



UCT - upper conf. bound  
for trees

$$\operatorname{argmax} \left[ \frac{w_i}{n_i} + c \sqrt{\frac{\ln(Z n_i)}{n_i}} \right]$$

