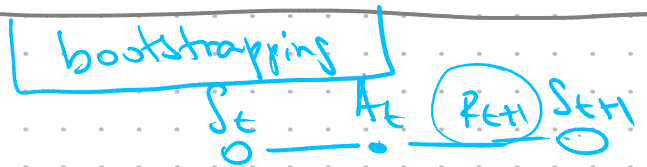


$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s]$$

$$V_{\pi'}, Q_{\pi'}$$

Temporal difference learning



$$V_{\pi}(s) = E_{\pi} [G_t | S_t = s] = E_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$G_t = R_{t+1} + \gamma G_{t+1} \approx R_{t+1} + \gamma V_{\pi}(S_{t+1})$$

TD(0) for estimation

Update ← Step · (New - Old)

- mit V_{π}

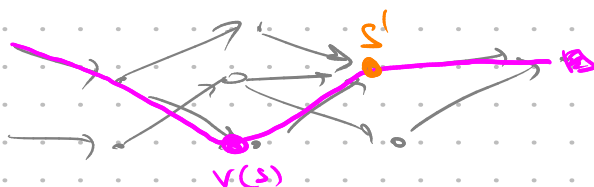
- repeat:

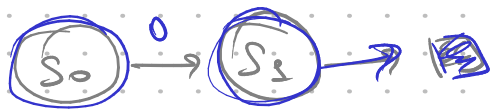
- erzeugen A_t & s_{t+1} nach π , receive R_{t+1}, S_{t+1}

- beobachten:

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

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





$$V(S_1) = \frac{1}{2}$$

$$S_1 \rightarrow \text{shaded} \quad R = 1, 0, 1, 1, 0$$

$$V^{MC}(S_0) = 0$$

$$S_0 \rightarrow S_1 \rightarrow \text{shaded} \quad R = 0$$

$$V^{TD}(S_0) = 0 + \gamma \cdot V^{TD}(S_1) = \frac{1}{2}$$

TD control

Sarsa (s, a, r, s', a')

On-policy TD control

- init
- repeat:

- bootstrap: A_t and S_t no π , nor R, S

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

gamma of Q, nonp.

$$\pi(A_t) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|A|} & \text{can go} \\ \frac{\epsilon}{|A|} & \text{argmax } Q \end{cases}$$

Expected Sarsa

$$+ \gamma \cdot \mathbb{E}_{\pi(a|S_{t+1})} [Q(S_{t+1}, a)] = \gamma \sum_a \pi(a|S_{t+1}) Q(S_{t+1}, a)$$

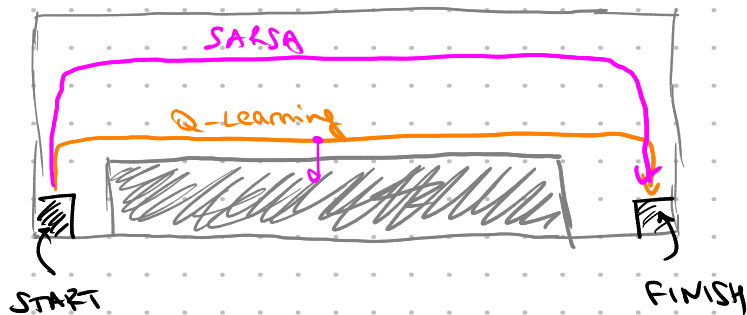
Off-policy TD control

Q-learning

1989

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

TD - Gammon

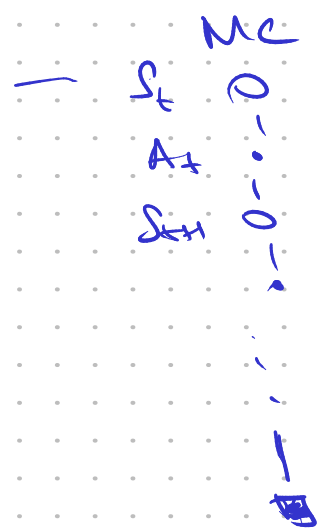
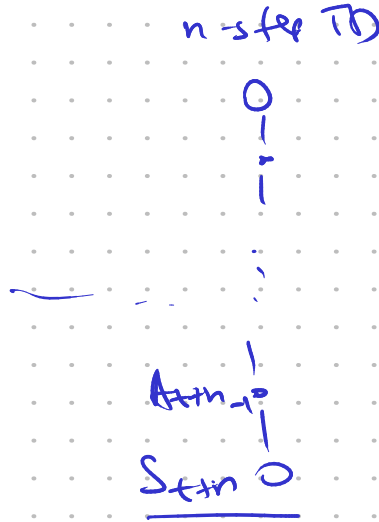
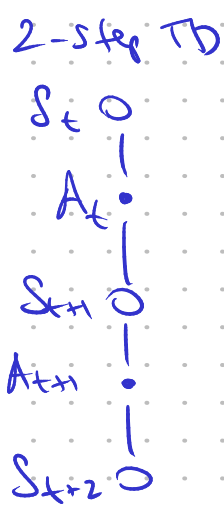
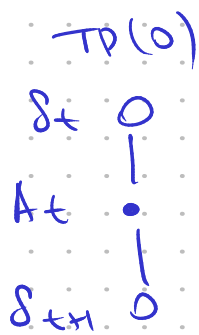


$$G_{t:t+1} = R_{t+1} + \gamma V(S_{t+1})$$

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

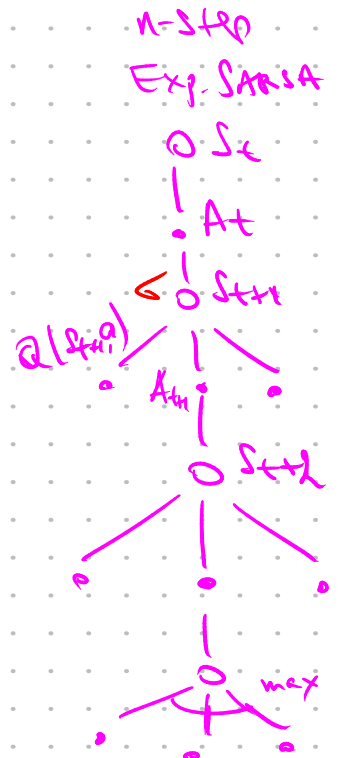
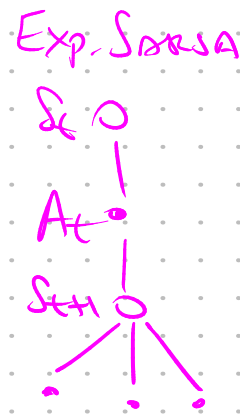
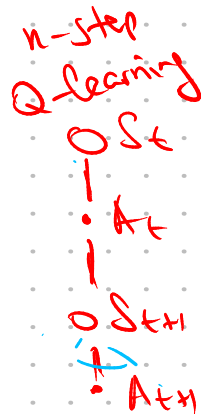
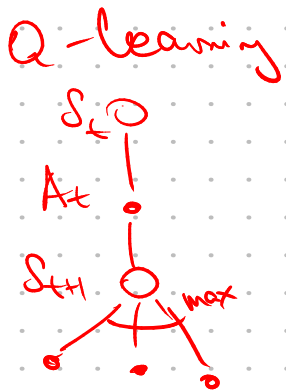
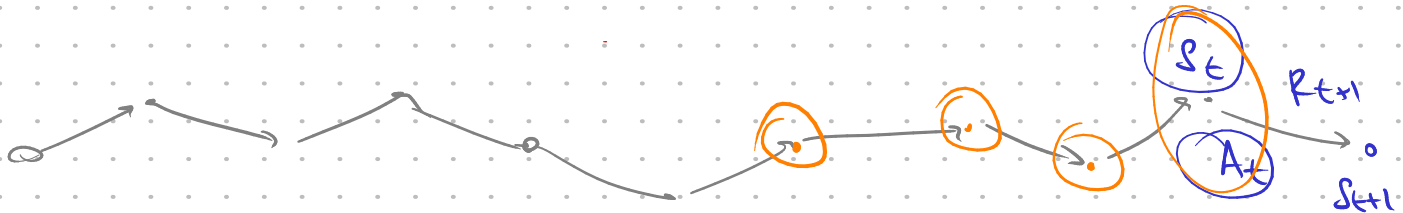
$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$$

$$G_{t:\infty} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n R_{t+n-1} + \dots = G_t$$



n-step SARSA:

$$\alpha (R_{t+n} + \gamma R_{t+n-1} + \dots + \gamma^{n-1} R_{t+1} + Q(S_{t+n}, A_{t+n}) - Q(S_t, A_t))$$



Sutton
Barto
1998/2018

$$Q(S_t, A_t) = Q(S_t, A_t) + \alpha \cdot (\hat{G} - Q(S_t, A_t))$$

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

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

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

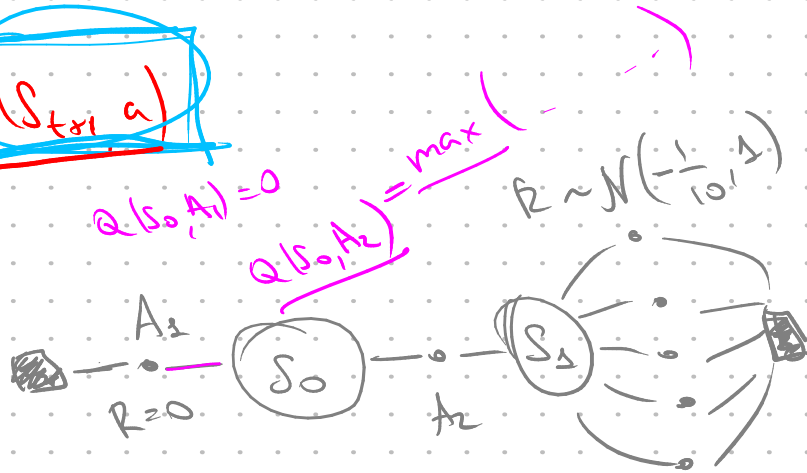
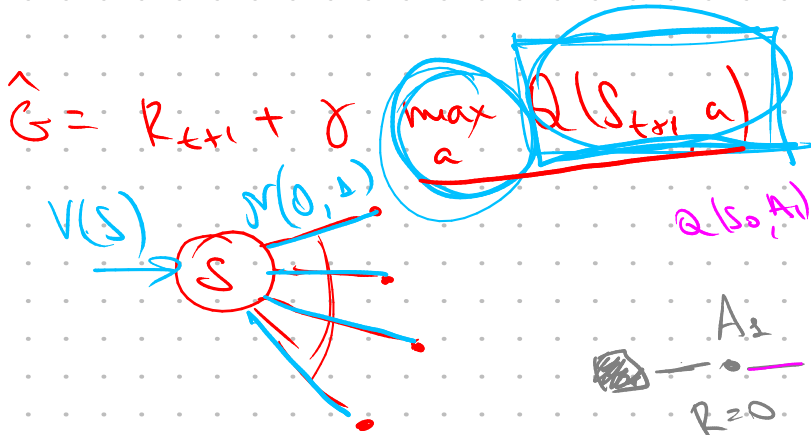
Double Q-learning

Winner's curse

Vickrey's auction



- 1: $\hat{x}_1 = x + N(0, \sigma^2)$
- 2: $\hat{x}_2 = x + N(0, \sigma^2)$
- 3: $\hat{x}_3 = x + N(0, \sigma^2)$



$$\hat{G}_{\text{double}} = R_{t+1} + \gamma Q_2(S_{t+1}, \operatorname{argmax}_a Q_1(S_{t+1}, a))$$

Double Q-learning:

- (s, a, r, s')

- Morkenka:

- $Q_1(s, a) := Q_1(s, a) + \alpha (r + \gamma Q_2(s', \operatorname{argmax}_a Q_1(s', a)) - Q_1(s, a))$

und

- $Q_2(s, a) := \dots$

