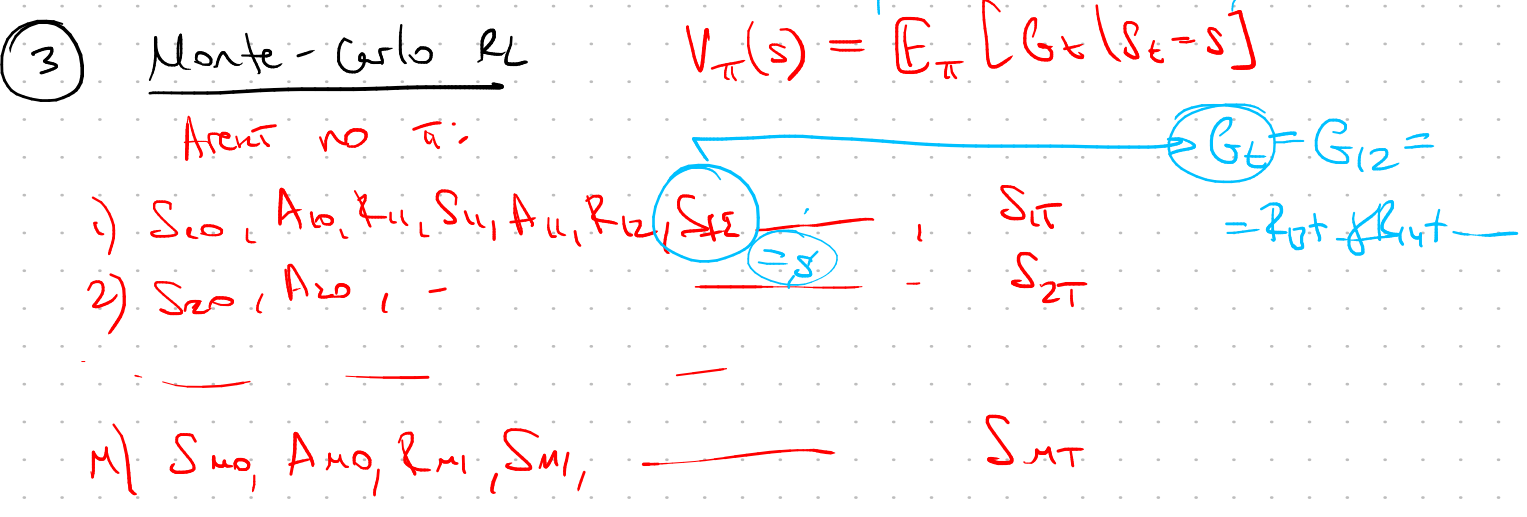
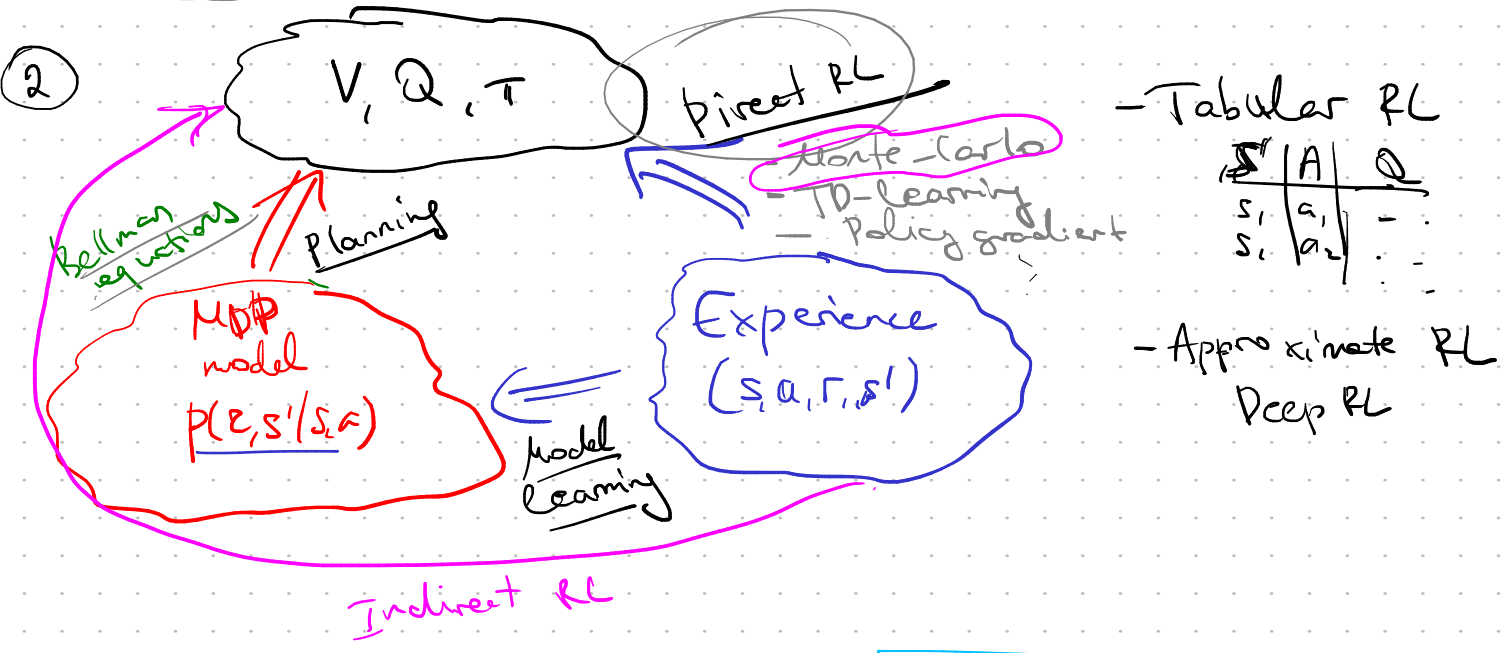
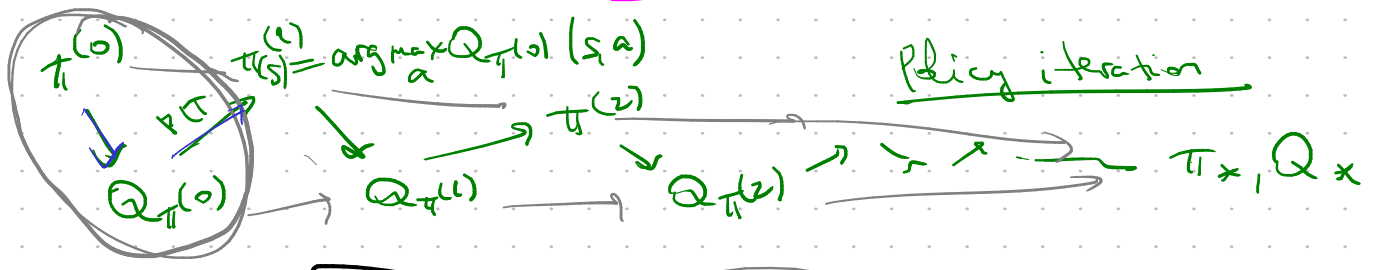
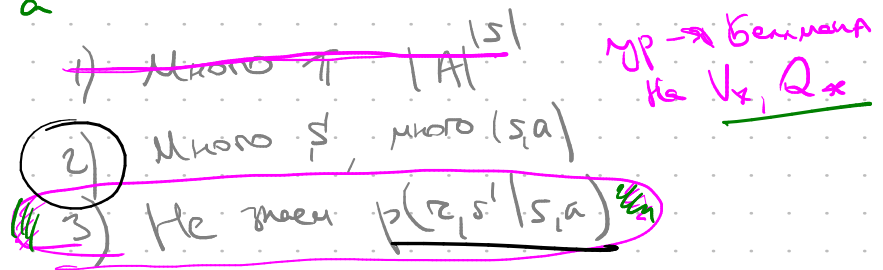


$V_\pi(s) = E_\pi[G_t | S_t = s]$ $V_*(s) = \max_\pi V_\pi(s)$
 $Q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]$ $Q_*(s, a) = \max_a Q_\pi(s, a)$
 $\pi_*(s) = \operatorname{argmax}_a Q_*(s, a)$

Bellman equations

Policy improvement theorem



Monte-Carlo estimation ($\pi \rightarrow V, Q$)

- init $V(s)$
- repeat:
 - generate w/ π $(S_t, A_t, R_{t+1}, S_{t+1}, \dots, S_T)$
 - $G = 0$
 - for $t = T-1, \dots, 0$
 - $G := \gamma G + R_{t+1}$
 - add G to the list $\text{Returns}(S_t)$ / $\text{Returns}(S_t, A_t)$
- $V_{\pi}(s) := \text{Avg}(\text{Returns}(s))$, $Q_{\pi}(s, a) := \text{Avg}(\text{Returns}(s, a))$

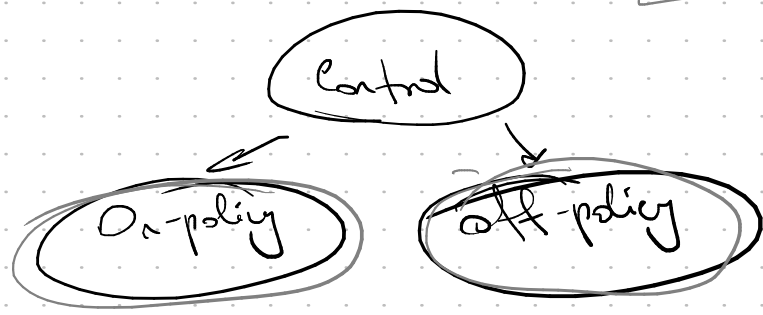
1) first visit MC
every-visit MC

2) Exploring starts

Monte-Carlo control with exploring starts

- init π, Q
- repeat:
 - generate episode $(S_t, A_t, R_{t+1}, S_{t+1}, \dots, S_m)$
 - for $t = T-1, \dots, 0$
 - add G_t to $\text{Returns}(S_t, A_t)$
 - $Q(S_t, A_t) := \text{Avg}(\text{Returns}(S_t, A_t))$
 - $\pi(S_t) := \text{argmax}_a Q(S_t, a)$

$(S_0, A_0) \sim \text{Unif}(s, a)$



On-policy MC control

$S_0 \sim \text{Unif}(\dots)$

$$\pi(S_t) := \begin{cases} \text{Unif}(a) & \epsilon \\ \text{argmax}_a Q(S_t, a) = a_x & 1 - \epsilon \end{cases}$$

$$\pi(a | S_t) = \frac{\epsilon}{|A|} \text{ unu } 1 - \epsilon + \frac{\epsilon}{|A|} \text{ go } a_x$$

$\forall s: Q_{\pi}(s, \pi'(s)) \geq \underline{V_{\pi}(s)} \quad \pi(\pi') - \epsilon - \text{soft}$

$$Q_{\pi}(s, \pi'(s)) = \sum_a \pi'(a|s) Q_{\pi}(s, a) =$$

$$= \frac{\epsilon}{|A|} \sum_a Q_{\pi}(s, a) + (1-\epsilon) \cdot \max_a Q_{\pi}(s, a)$$

$$V_{\pi}(s) = \sum_a \pi(a|s) Q_{\pi}(s, a) = \frac{\epsilon}{|A|} \sum_a Q_{\pi}(s, a) + \sum_a \left(\pi(a|s) - \frac{\epsilon}{|A|} \right) Q_{\pi}(s, a)$$

$$\sum_a \frac{\pi(a|s) - \frac{\epsilon}{|A|}}{1-\epsilon} = 1$$

$$(1-\epsilon) \cdot \sum_a \frac{\pi(a|s) - \frac{\epsilon}{|A|}}{1-\epsilon} Q_{\pi}(s, a) \leq \max_a Q_{\pi}(s, a)$$

$$\leq Q_{\pi}(s, \pi'(s))$$

④ Off-policy MC control

- generate episodes with π' , e.g. ϵ -soft
- learn Q_{π} for $\pi \neq \pi'$, e.g. greedy

$$Q_{\pi} = \mathbb{E}_{\pi} [G_{t+1} | s_t]$$

Importance sampling

$$\mathbb{E}_{p(\bar{x})} [f(\bar{x})] = \int p(\bar{x}) f(\bar{x}) d\bar{x} =$$

$$= \int f(\bar{x}) \frac{p(\bar{x})}{q(\bar{x})} q(\bar{x}) d\bar{x} = \mathbb{E}_q \left[f \left(\frac{p}{q} \right) \right]$$

Importance weights

$$p(\text{Traj} | \pi) = p(A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}, \dots | S_t = s, \pi) =$$

$$= \underbrace{p(A_t | S_t)}_{\pi} \underbrace{p(R_{t+1}, S_{t+1} | S_t, A_t)}_{p(R|S'(S, a))} \underbrace{p(A_{t+1} | S_{t+1})}_{\pi} \underbrace{p(R_{t+2}, S_{t+2} | S_{t+1}, A_{t+1})}_{\dots}$$

$$\mathbb{E}_{\pi} [G_t | S_t] = \mathbb{E}_{\pi'} \left[G_t \cdot \frac{p(\text{Traj} | \pi)}{p(\text{Traj} | \pi')} \right] =$$

$$= \mathbb{E}_{\pi'} \left[G_t \cdot \frac{\pi(A_t | S_t) \cancel{p(R_{t+1}, S_{t+1} | S_t, A_t)} \pi(A_{t+1} | S_{t+1}) \cancel{p(\dots)}}{\pi'(A_t | S_t) \cancel{p(R_{t+1}, S_{t+1} | S_t, A_t)} \pi'(A_{t+1} | S_{t+1}) \cancel{p(\dots)}} \right]$$

$$V_{\pi}(s_k) = E_{\pi'} [G_t + \gamma \sum_{k=t}^{T-1} \frac{\pi(A_k|s_k)}{\pi'(A_k|s_k)}] = w_t^{\pi, \pi'}$$

Off-policy MC control

- init $\pi'(a|s) = 0 \Rightarrow \pi(a|s) = 0$

- repeat:

- generate w/π' : $S_0, A_0, R_0, S_1, A_1, \dots$

- $G := 0, w := 1$

- for $t = T-1, \dots, 0$:

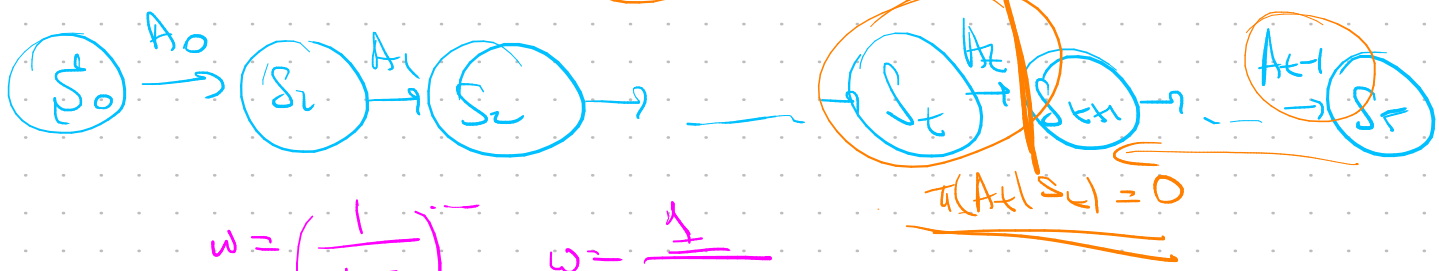
- $G := \gamma G + R_t$

- $w := w \cdot \frac{\pi(A_t|S_t)}{\pi'(A_t|S_t)}$

ecm $\pi(A_t|S_t) = 0, \pi' > 0$
 $w_t^{\pi, \pi'} = 0 \forall t' \leq t$

- add $G \cdot w$ to returns(S_t, A_t)

- $\pi' := \epsilon$ -soft(Q_{π}), $\pi'(a|S_t) = \begin{cases} \epsilon/|A|, \epsilon \\ 1-\epsilon, a = \text{argmax } Q_{\pi} \end{cases}$



$$w = \left(\frac{1}{1-\epsilon} \right)^t$$

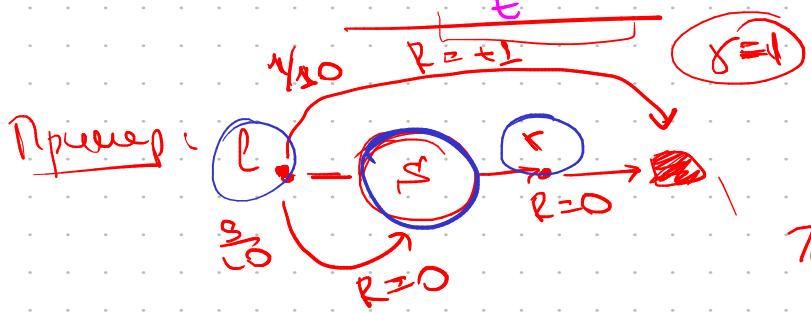
$$w = \frac{1}{\prod_{t'=t} \pi'(A_{t'}|s_{t'})}$$

1) Ordinary Emp. s.

$$V_{\pi}(s) := \text{Avg}(\text{Returns}(s)) = \frac{1}{|Ret|} \sum_{t: S_t=s} G_t \cdot w_t^{\pi, \pi'}$$

2) Weighted Emp. s.:

$$V_{\pi}(s) := \frac{\sum_t G_t \cdot w_t^{\pi, \pi'}}{\sum_t w_t^{\pi, \pi'}}$$



$Q(s, \text{right}) = 0$
 $Q(s, \text{left}) = 1$

$\pi = \text{left}, \pi' = \left(\frac{1}{2}, \frac{1}{2} \right)$

$$E_{\pi} [G_t \cdot w_t^{\pi, \pi}] = 1 \quad \text{Var} [\dots] = ?$$

$$\text{Var} [x] = E[x^2] - \underbrace{E[x]^2}_{= 1} = E[x^2] - 1$$

$$E_{\pi} \left[G_0 \cdot \prod_{t=0}^{T-1} \frac{\pi(A_t | S_t)}{\pi'(A_t | S_t)} \right] = \sum_{\text{Traj}} p(\text{Traj} | \pi) \cdot (w)^2 \cdot G^2$$

$$= \sum_{\text{Traj}: G=1} (w_0^{\pi, \pi})^2 p(\text{Traj} | \pi) =$$

$$= \frac{1}{2} \cdot \frac{1}{10} \cdot \left(\frac{1}{4}\right)^2 + \underbrace{\left(\frac{1}{2}\right)^2 \frac{9}{10} \cdot \frac{1}{10}}_{l, l+1} \cdot \left(\frac{1}{2} - \frac{1}{2}\right)^2 + \left(\frac{1}{2}\right) \cdot \left(\frac{9}{10}\right) \cdot \frac{1}{10} \left(\frac{1}{2} - \frac{1}{2}\right)^2$$

$$= \sum_{t=0}^{\infty} \frac{1}{10} \cdot \left(\frac{9}{10}\right)^t \cdot \left(\frac{1}{2}\right)^{t+1} \cdot 2^{2(t+1)} = \frac{2}{10} \cdot \sum_t \left(\frac{9}{5} \cdot \frac{1}{2} \cdot 4\right)^t =$$

$$= \frac{1}{5} \cdot \sum_t \left(\frac{9}{5}\right)^t$$



$$R = -1$$

$$\bar{P} = (P_s - P_e)$$

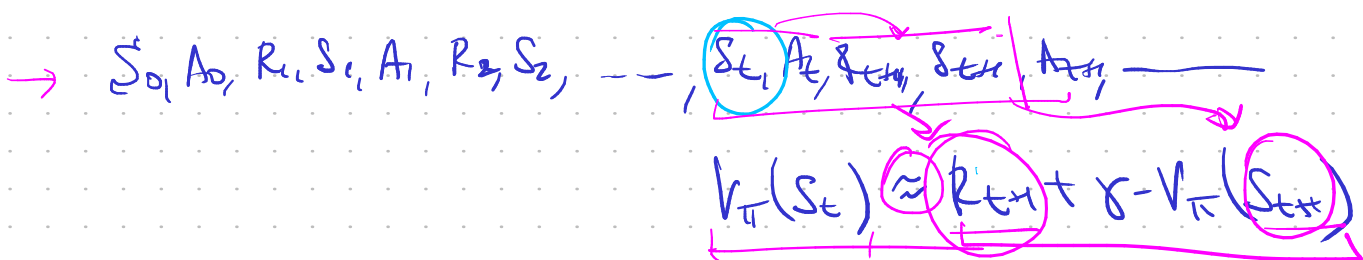
-1 → Win $P_e \rightarrow P_s$

⑤ TD-learning (temporal difference)

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

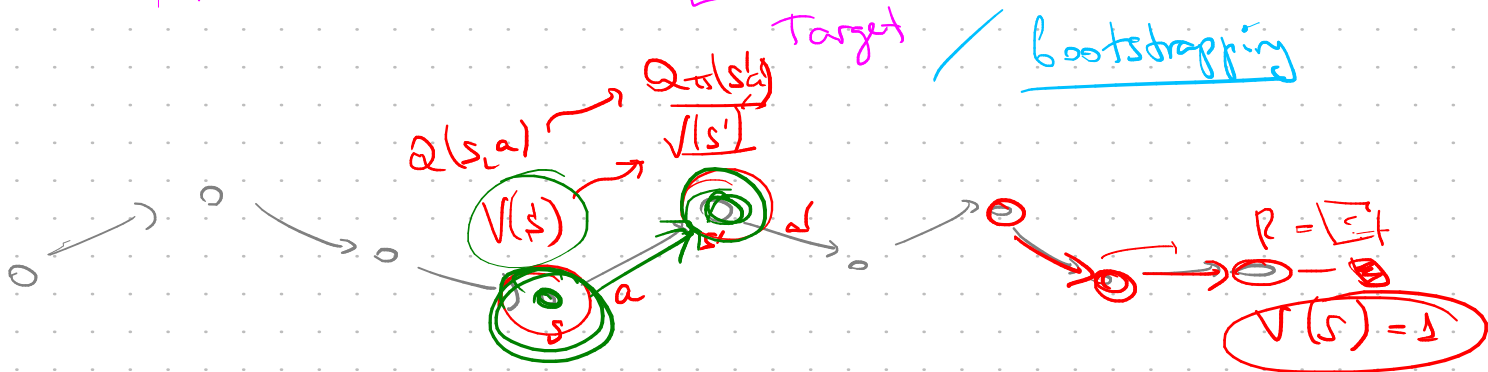
$$= \underbrace{E_{\pi} [R_{t+1}]}_{\approx R_{t+1}} + \gamma \cdot \underbrace{E_{\pi} [V_{\pi}(S_{t+1})]}_{\approx G_{t+1}}$$





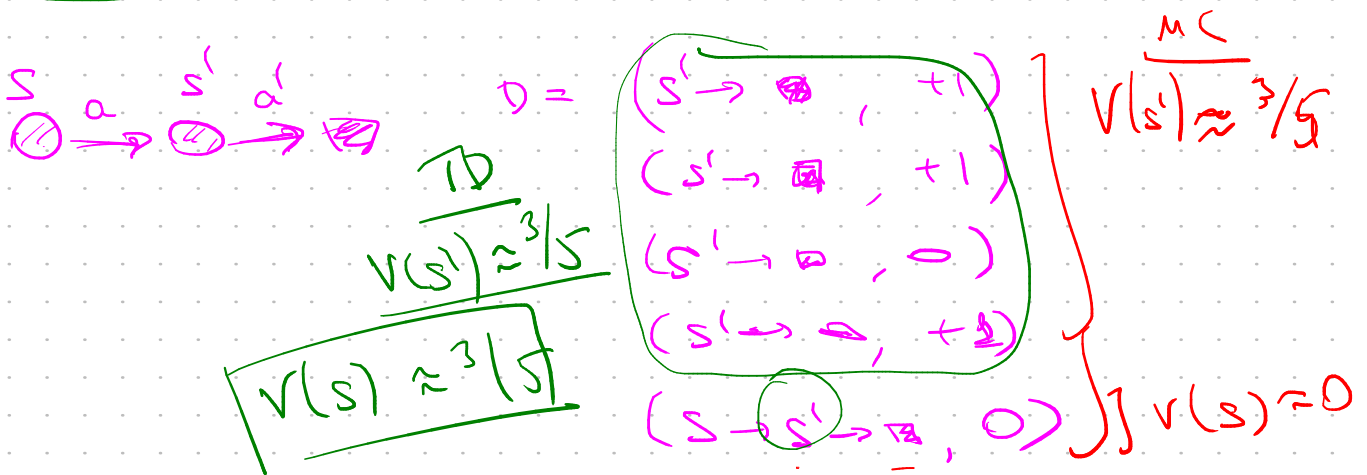
Experience tuple: (s, a, r, s')

$V_{\pi}(s) := V_{\pi}(s) + \alpha (\underbrace{r + \gamma V_{\pi}(s')}_{\text{Target}} - V_{\pi}(s))$



TD estimation

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



Backup diagrams

Bellman eq: $V(s) = \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) (r + \gamma V(s'))$

