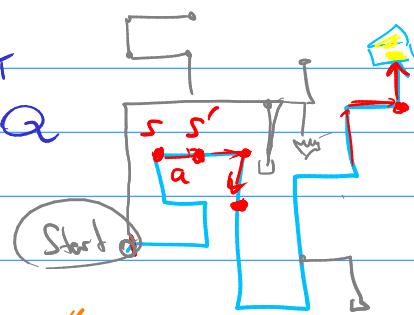
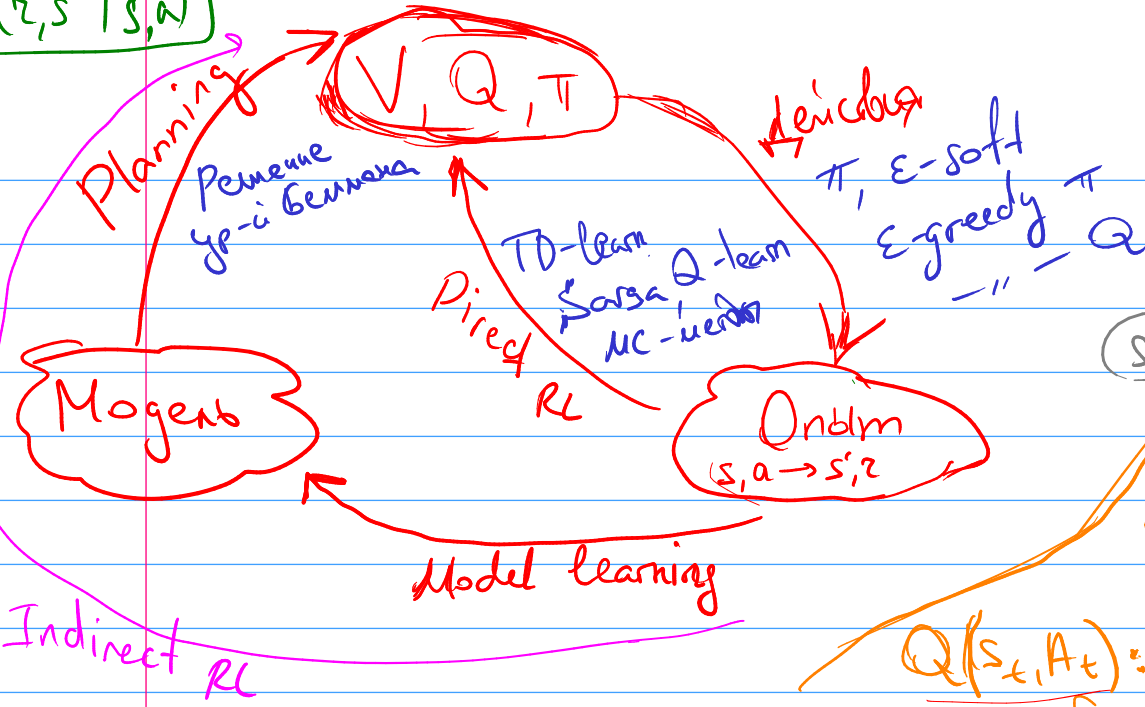


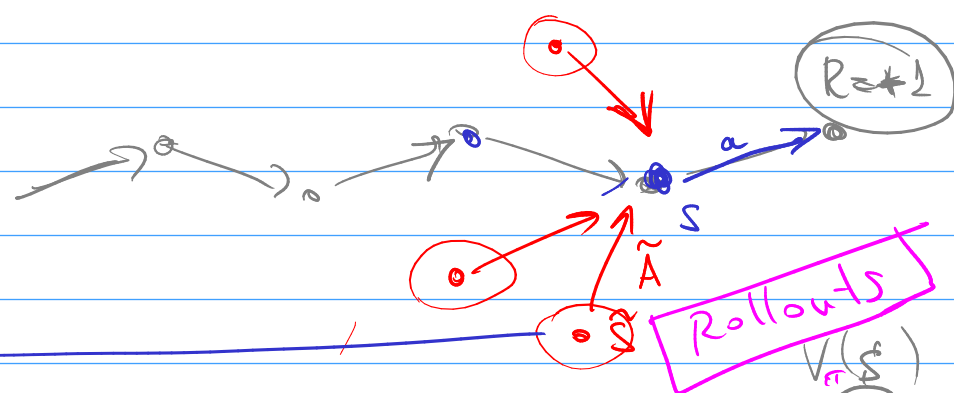
$$p(z, s' | s, a)$$

$$V(s), Q(s, a)$$



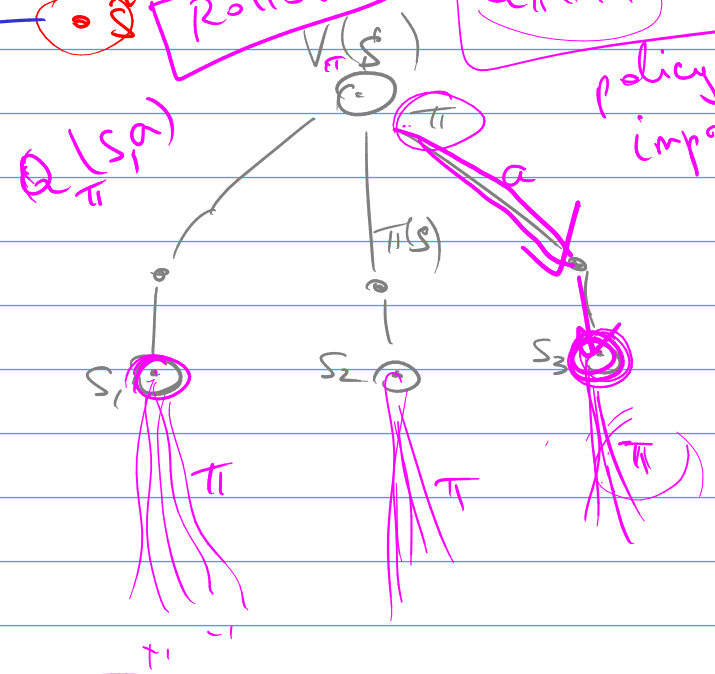
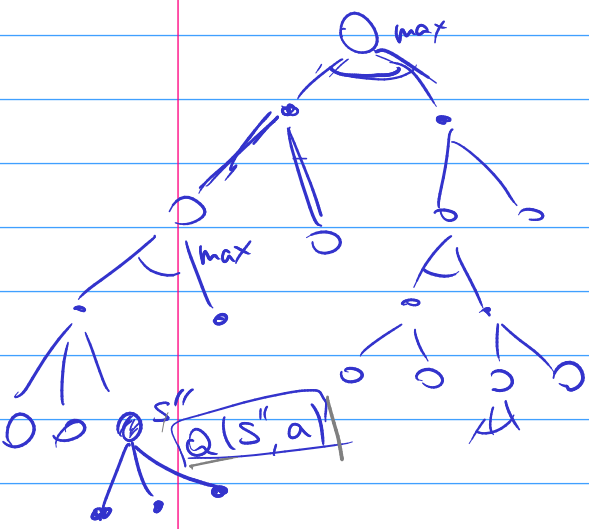
$$(s, a) \rightarrow (s', z)$$

$$Q(s_t, A_t) := Q(s_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, A_t)]$$



$$Q_\pi(s, a) \approx V_\pi(s)$$

policy improvement



MCTS - Monte-Carlo Tree Search

n-step methods

$$Q(S_t, A_t) = Q(S_t, a) + \alpha (\text{TARGET} - Q(S_t, a))$$

- t
- t+1
- t+2
- t+n

MC

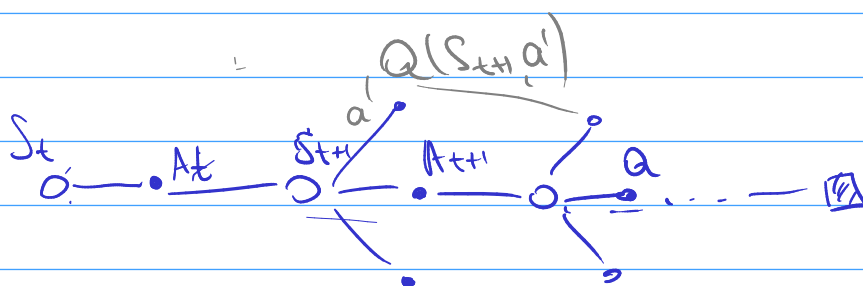
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{T-t-1} R_T$$

TD:

$$G_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$$

2-step TD:

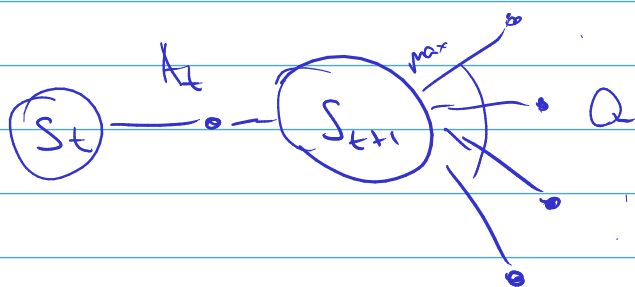
$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 Q(S_{t+2}, A_{t+2})$$



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

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

$$+ \gamma \pi(A_{t+2}|S_{t+2}) (\dots)$$



$$S \rightarrow \boxed{NN} \rightarrow V(s)$$

$$\begin{matrix} S \\ a \end{matrix} \rightarrow \boxed{NN} \rightarrow Q(s, a)$$

$$S \rightarrow \boxed{NN} \rightarrow \begin{matrix} Q(s, a) \\ Q(s, a_t) \end{matrix} \leftarrow$$