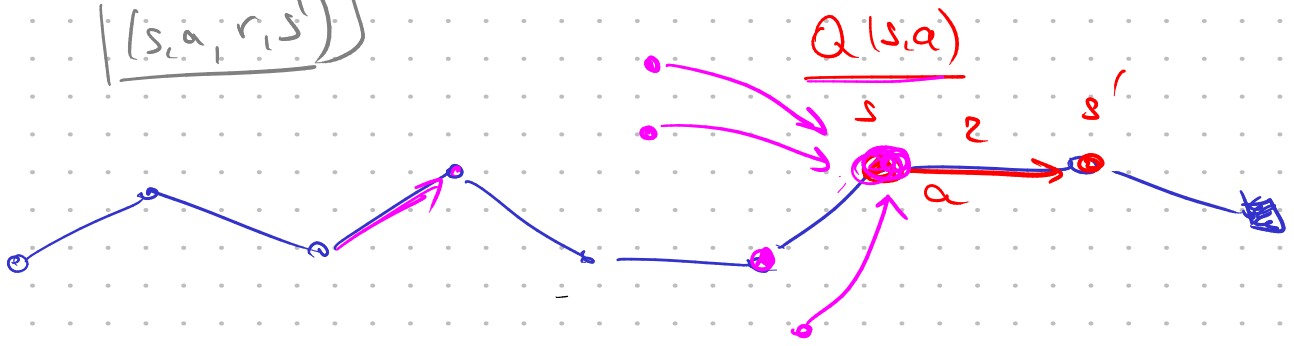


$$Q(s,a) := Q(s,a) + \alpha (r + \gamma \cdot \max_{a'} Q(s',a') - Q(s,a))$$

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

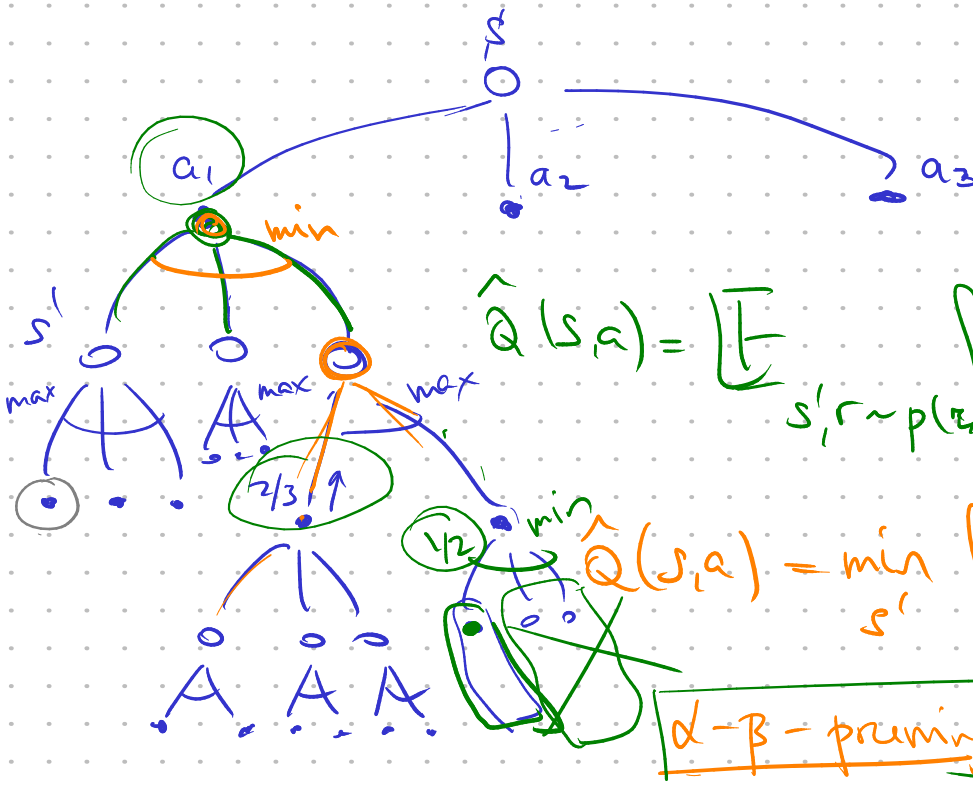


$V(s)$

$Q(s,a)$

Decision-time
planning

$$\pi(s) = \operatorname{argmax}_{a'} Q(s,a')$$

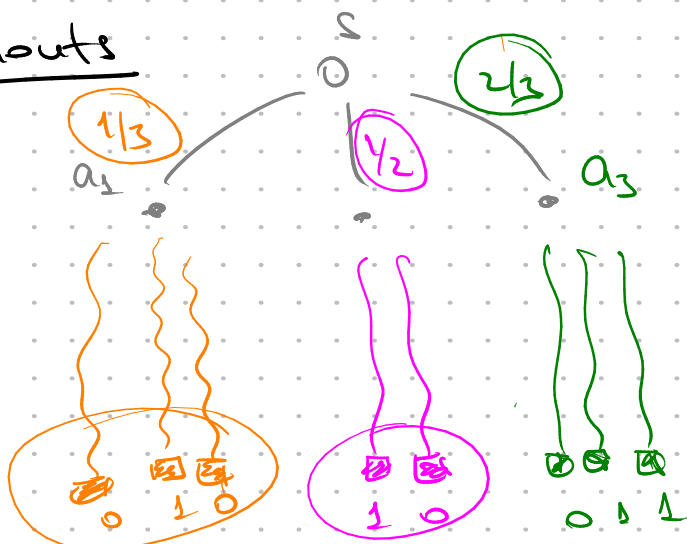


$$\hat{Q}(s,a) = \mathbb{E}_{s', r \sim p(s', r | s, a)} [r + \gamma \cdot \max_{a'} Q(s', a')]$$

$$\hat{Q}(s,a) = \min_{s'} [r + \gamma \cdot \max_{a'} Q(s', a')]$$

alpha-beta-pruning

Rollouts



$$\hat{Q}(s, a_3) = \alpha Q(s, a_3) + (1 - \alpha) \cdot \hat{R}$$

TD-Gammon

Tesauro

MCTS - Monte Carlo Tree Search

