"Markov Decision Process"

$S$ : all possible states

$A$ : " " actions

$R$ : reward distribution given $(s, a)$

$P$ : transition prob. to $S_{t+1}$ given $(s, a)$

$\gamma$ : discount factor

$$S \xrightarrow{\pi} A$$

objective:
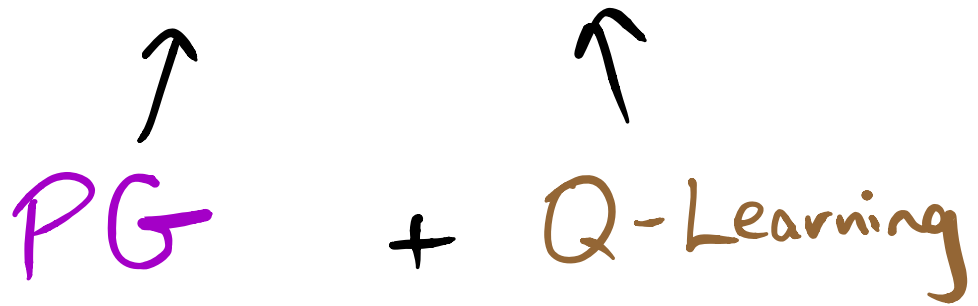
find $\pi^* = \max \left( \sum_{t>0} \gamma^t r^t \right)$

value function: $V^\pi(s)$

Q-value function: $Q^\pi(s, a)$

$Q^*(s, a) \approx Q(s, a, \theta)$

# Actor - Critic

$$\text{PG} + \text{Q-Learning}$$