

Q-LEARNING

- $Q^{\pi}(s,a) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k r_k \mid s, a \right]$
- $Q^{\pi}(s,a) = \mathbb{E}_{s',a',r}[R(s,a) + \gamma Q^{\pi}(s',a') | s,a]$
- $\bullet \ \overline{Q^*(s,a)} = \mathbb{E}_{s',r} \left[R(s,a) + \max_{a'} \overline{Q^*(s',a')} \, | s,a \right]$
- $\hat{Q}(s,a) \leftarrow r(s,a) + \gamma \max_{a'} \hat{Q}(s',a')$
- $\hat{Q}(s,a) \leftarrow \hat{Q}(s,a) + \alpha [r(s,a) + \gamma \max_{a'} \hat{Q}(s',a') \hat{Q}(s,a)]$

Q-LEARNING: OFF POLICY

Q-learning (off-policy TD control) for estimating $\pi \approx \pi_*$

```
Algorithm parameters: step size \alpha \in (0,1], small \varepsilon > 0

Initialize Q(s,a), for all s \in \mathbb{S}^+, a \in \mathcal{A}(s), arbitrarily except that Q(terminal, \cdot) = 0

Loop for each episode:

Initialize S

Loop for each step of episode:

Choose A from S using policy derived from Q (e.g., \varepsilon-greedy) Exploration

Take action A, observe R, S'

Q(S,A) \leftarrow Q(S,A) + \alpha \left[R + \gamma \max_a Q(S',a) - Q(S,A)\right] Exploitation

S \leftarrow S'

until S is terminal
```

Q-LEARNING -> DQN

Pros

- Same principal but use function approximation
- Large scale problems, huge (Continuous) state space
- More flexible data(e.g. Images, signals)
- Stability boost: Replay buffer & Target network
- Efficiency boost: Output array of Qs

Issues

- Limited to discrete action space
- Bad for large action space

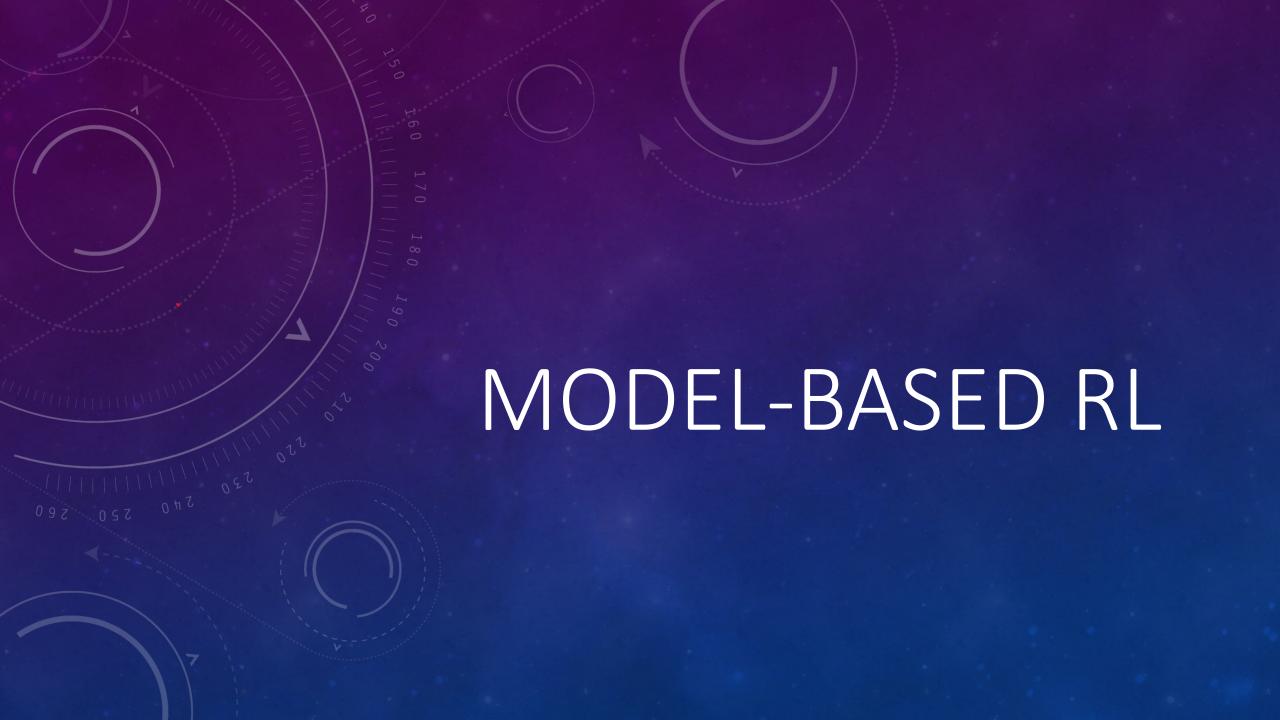
DQN IMPLEMENTATION

$\overline{\mathbf{Algorithm}}$ 4 DQN

```
1: procedure DQN
           Initialize network Q_{\omega} and Q_{\text{target}} as a clone of Q_{\omega}
          Initialize replay buffer R and burn in with trajectories followed by random policy
 3:
          repeat for E training episodes:
                 Initialise S_0
 5:
                 for t = 0, 1, ..., T - 1:
 6:
                        a_t = \begin{cases} \arg\max_a Q_{\omega}(s_t, a) & \text{with probability } 1 - \epsilon \\ \text{Random action} & \text{otherwise} \end{cases}
 7:
                         Take a_t and observe r_t, s_{t+1}
                         Store (s_t, a_t, r_t, s_{t+1}) in R
 9:
                         Sample minibatch of (s_i, a_i, r_i, s_{i+1}) with size N from R
10:
                       y_i = \begin{cases} r_i & s_{i+1} \text{ is terminal} \\ r_i + \gamma \max_a Q_{\text{target}}(s_i, a) & \text{otherwise} \end{cases}
11:
                        L(\omega) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - Q_{\omega}(s_i, a_i))^2 -
12:
                         Update\ Q_{\omega}\ using\ {	t Adam}\ (
abla_{\omega}L(\omega))
13:
                        Replace Q_{\text{target}} with current Q_{\theta} if t \mod 50 = 0
14:
15: end procedure
```

Loss reduced from 2D data

ONLY update on ONE ACTION, NO CHANGE to other Qs



IDEA

- Learns an environment/dynamics model/network
- Function Approximation: $P_{\varphi}: S \times A \rightarrow R \times S'$
- Model rollout: $S_t \to A_t | \pi_{\theta}(S_t) \to R_t | P_{\varphi}(S_t, A_t) \to S_{t+1} | P_{\varphi}(S_t, A_t) \to \cdots$
- (Probabilistic) Ensemble networks

MBRL VS MODEL FREE: INTUITION

- Psychology: Unconscious vs Conscious
- Biology: Habitual vs Goal-directed





MBRL VS MODEL FREE

- More sample efficient(Short term)
- Compounding error(Stochasticity & Environment variety)
- Challenge: State representation & Use of prior domain knowledge

 Legged robot application: MBRL to stand and Model free RL to walk(Surrounding obstacles)