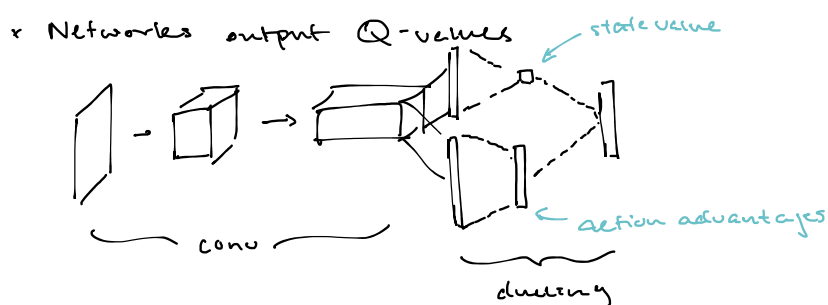


DUELING NETWORK ARCHITECTURES FOR DEEP REINFORCEMENT LEARNING

► model-free RL

- * Factor state value function V & state dependent actor advantage function A into two separate streams.



* $V(s), A(s,a)$

Advantage updating converged faster than Q -learning

2. Background

- * Agent gets input (M frames) of $\mathbf{z}_t = (x_{t-M+1}, \dots, x_t) \in \mathcal{Z}$

- * maximize expected discounted return.

$$E[R_t = \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{\tau}]$$

- * $Q^{\pi}(s,a) = E[R_t | s_t=s, a_t=a, \pi]$

$$V^{\pi}(s) = E_{a \sim \pi(s)}[Q^{\pi}(s,a)]$$

$$Q^{\pi}(s,a) = E_{s'}[r + \gamma E_{a' \sim \pi(s')}[Q^{\pi}(s',a')] | s,a,\pi]$$

For optimal Q value, Q^*

(w/ deterministic policy)

$$a = \operatorname{argmax}_{a' \in \mathcal{A}} Q^*(s,a') \Rightarrow E_{a' \sim \pi(s')} [Q^*(s',a')] = \max_{a' \in \mathcal{A}} Q^*(s',a')$$

$$Q^*(s,a) = E_{s'}[r + \gamma \max_{a'} Q^*(s',a') | s,a]$$

- * Advantage function

$$A^{\pi}(s,a) = Q^{\pi}(s,a) - V^{\pi}(s)$$

$$E_{a \sim \pi(s)}[A^{\pi}(s,a)] = 0 \quad \text{PROVE THIS!}$$

a depends on s.

2.1 DQN

- * NN: $Q(s,a;\theta)$.

- * At iteration i , loss function is

$$L_i(\theta_i) = E_{s,a,r,s'}[(y_i^{\text{DQN}} - Q(s,a;\theta_i))^2]$$

$$y_i^{\text{DQN}} = r + \gamma \max_{a'} Q(s',a';\theta^-)$$

target is the value of the best action/state in the target network (frozen weights).

$$\nabla_{\theta_i} L_i(\theta_i) = E_{s,a,r,s'}[(y_i^{\text{DQN}} - Q(s,a;\theta_i)) \nabla_{\theta_i} Q(s,a;\theta_i)]$$

$(s,a,r,s') \sim \mathcal{U}(\mathcal{D})$ reduces correlation among samples.
 dataset of experiences.

2.2 Double DQN [van Hasselt et al (2015)]

- * Q-learning & DQN, 'max' operator uses same values to both select & evaluate an action.

$$y_i^{\text{DDQN}} = r + \gamma Q(s', \operatorname{argmax}_{a'} Q(s,a;\theta_i); \theta^-)$$

y_i^{DDQN} selected the action \hat{a} evaluated it — bias!

2.3 Prioritized Replay [Schuul 2016]

- * increase replay probability of experience tuples that have high expected learning progress (proxy: absolute TD-error).

3 Dueling Network Architecture

- * evaluating value of action choice doesn't matter for many states.

* $Q^{\pi}(s,a) = V^{\pi}(s) + A^{\pi}(s,a)$

$$V^{\pi}(s) = E_{a \sim \pi(s)}[Q^{\pi}(s,a)]$$

$$\Rightarrow E[A^{\pi}(s,a)] = 0$$

- * $V(s;\theta,\beta) \rightarrow \beta$ is fully connected

- * $A(s,a;\theta,\alpha) \rightarrow \alpha$ is fully connected

* $Q(s,a;\theta,\alpha,\beta) =$

$$V(s;\theta,\beta) + A(s,a;\theta,\alpha)$$

bad b/c we can't identify V & A uniquely ($V+\delta + A-\delta$) $\forall \delta$.

remedy; subtract

$$\left[\max_{a' \in \mathcal{A}} A(s,a';\theta,\alpha) \right] \text{ (chosen action)}$$

4. Experiments

4.1 Policy Evaluation

if we subtract $Q(s,a;\theta,\alpha,\beta) = V(s;\theta,\beta)$ instead, Q will be off-target, but more stable.

