Variational Deep Q-Networks in Edward

Harri Bell-Thomas

R244: Open Source Project Presentation

19/11/2019

・ロト・日本・ヨト・ヨト・日・ つへぐ

Q-Learning

Q-Learning is *model-free* reinforcement learning.

 ${\mathcal Q}$ is the action-value function defining the reward used for reinforcement — this is learned.

Conceptually,

$$\mathcal{Q}^{\pi}(s, a) = \mathbb{E}_{a_t \sim \pi(\cdot | s_t)} \Big[\sum_{t=0}^{\infty} r_t \gamma^t \, | \, s_0 = s, a_0 = a \Big]$$

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

The value of Q^{π} at a certain point in time, t, in terms of the payoff from an initial choice, a_t , and the value of the remaining decision problem that results after that choice.

$$J(\pi) = \mathbb{E}_{\alpha} \Big[\big(\mathcal{Q}^{\pi}(s_t, a_t) - \max_{a} \mathbb{E}[r_t + \gamma \mathcal{Q}^{\pi}(s_{t+1}, a)] \big)^2 \Big]$$

$$\alpha \to s_0 \sim \rho, \, a_t \sim \pi(\cdot | s_t)$$

・ロト・日本・モト・モー ショー ショー

Deep Q-Networks

Briefly

Approximate the action-value function $Q^{\pi}(s, a)$ with a neural network $Q_{\theta}(s, a)$. The (greedy) policy represented by this is π_{θ} .

Discretise the expectation using K sample trajectories, each with period T. Use this to approximate $J(\theta)$.

$$\tilde{J}(\theta) = \frac{1}{K} \frac{1}{T} \sum_{i=1}^{K} \sum_{t=0}^{T-1} (\mathcal{Q}_{\theta}^{(i)}(s_{t}^{(i)}, a_{t}^{(i)})) - \max_{a} \left[r_{t} + \gamma \mathcal{Q}_{\theta}^{(i)}(s_{t+1}^{(i)}, a) \right])^{2}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○ ● ●

Main Concepts:

1. Try to solve an optimisation problem over a class of tractable distributions, q, parameterised by ϕ , in order to find the one most similar to p.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○ ● ●

- 2. $\phi \leftarrow \min_{\phi} \mathbb{KL}(q_{\phi}(\theta) \parallel p(\theta|D))$
- 3. Approximate this using gradient descent.

Variational Deep Q-Networks

Idea: For efficient exploration we need $q_{\phi}(\theta)$ to be dispersed near even coverage of the parameter space. Encourage this by adding an entropy bonus to the objective.

$$\mathbb{E}_{\theta \sim q_{\phi}(\theta)} \Big[\big(\mathcal{Q}_{\theta}(s_{j}, a_{j}) - \max_{a'} \mathbb{E}[r_{j} + \gamma \mathcal{Q}_{\theta}(s'_{j}, a')] \big)^{2} \Big] - \lambda \mathbb{H}(q_{\phi}(\theta))$$

Assigning systematic randomness to Q enables efficient exploration of the policy space. Further, encouraging high entropy over parameter distribution prevent premature convergence.

tl;dr Higher chance of finding maximal rewards in a faster time than standard DQNs.

Algorithm

Algorithm 1 Variational DQN

- 1: INPUT: improper uniform prior $p(\theta)$; target parameter update period τ ; learning rate α ; generative model variance σ^2
- 2: INITIALIZE: parameters ϕ , ϕ^- ; replay buffer $R \leftarrow \{\}$; step counter counter $\leftarrow 0$
- 3: for e = 1, 2, 3...E do
- 4: while episode not terminated do
- 5: $counter \leftarrow counter + 1$
- 6: Sample $\theta \sim q_{\theta}(\phi)$
- 7: In state s_t , choose $a_t = \arg \max_a Q_{\theta}(s_t, a)$, get transition s_{t+1} and reward r_t
- 8: Save experience tuple $\{s_t, a_t, r_t, s_{t+1}\}$ to buffer R
- 9: Sample N parameters $\theta_i^- \sim q_\theta(\phi^-)$ and sample N tuples $D = \{s_j, a_j, t_j, s'_i\}$ from R
- 10: Compute target $d_j = r_j + \max_{a'} Q_{\theta_j^-}(s'_j, a')$ for *j*th tuple in D
- 11: Take gradient $\Delta \phi$ of the KL divergence in (8)

12:
$$\phi \leftarrow \phi - \alpha \Delta \phi$$

- 13: **if** counter mod $\tau = 0$ **then**
- 14: Update target parameter $\phi^- \leftarrow \phi$
- 15: end if
- 16: end while
- 17: end for

Figure: VDQN Pseudocode.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

Aim / Goals

Workplan



Questions?

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへで