

- Policy evaluation: Approximate policy evaluation $Q(s, a, \pi) \approx Q^{\pi}(s, a)$
- Policy improvement: ϵ -greedy policy improvement

DQN

DQN: Bringing Deep Learning into RL I

In principle it sounds easy: use convolutional neural networks to link screen shots to values.

Given our TD Q-learning update:

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a') - Q(s_t, a)] \quad (1)$$

we want to learn $Q(s_t, a; w)$ as a parametrised function (a neural network) with parameters w .

We can define the TD error as our learning target that we want to reduce to zero.

$$TD\ error(w) = r_{t+1} + \gamma \max_{a_{t+1}} Q(s_{t+1}, a'; w) - Q(s_t, a; w) \quad (2)$$

Taking the gradient of E wrt w we obtain:

$$\Delta w = [\mathbf{r} + \gamma \max Q(s_{t+1}, a_{t+1}; w) - Q(s_t, a; w)] \nabla_w Q(s_t, a; w) \quad (3)$$

However, the ATARI playing problem required more engineering to solve properly

- Experience Replay
- Target Network
- Clipping of Rewards
- Skipping of Frames

DQN: Experience Replay I

The CNN is easily overfitting the latest experienced episodes and the network becomes worse at dealing with different experiences of the game world:

- Complication 1: Inefficient use of interactive experience

- Training deep neural networks requires many samples, each has its own transition
- But now each state transition is used only once, then discarded
- So we need to constantly revisit the same state transitions to train

This is inefficient and slow.

- Complication 2: Experience distribution

- Interactive learning produces training samples that are highly correlated (agents recent actions reflect recent policy)
- The probability distribution of the input data changes

- The network forgets transitions that were in the not so recent past (overwriting past memory, in neuroscience known as "extinction")

- Moreover, because networks are monolithic, changes in one part of the network can have side-effects on other parts of the network [Why is this not a problem for table-based RL?]

Solution: Experience Replay² (alternative use Hierarchical RL)

- Experiences (traces) are stored and replayed in mini-batches to train the neural networks on more than just the last episode.

- Instead of running Q-learning on each state-action pair as they occur during simulation or actual experience, the experience replay buffer stores the traces sampled.

- Mini-batches are only feasible when running multiple passes with the same data is somewhat stable with respect to the samples. I.e. the transitions should show low variance/entropy for the next immediate outcomes (reward, next state) given the same state, action pair.

- Replaying past data is also a more efficient use of previous experience, by learning (training the neural network) with it multiple times. This is key when gaining real-world experience is costly (e.g. simulation time) as the Q-learning updates are incremental and do not converge quickly.

The learning phase is separated from gaining experience, and based on random samples from the mini-batch table. DQN interleaves the two processes – acting and learning. Improving the policy will lead to a different behaviour (and different state action pairs), resulting too large replay buffers getting stale.

In summary

- Reduction of correlation between experiences in updating DNN – mini-batches effectively make the samples more independent and identically distributed.

- Increased learning speed with mini-batches due to reuse/multi-use of experience

- Reusing/Replacing past transitions to avoid catastrophic forgetting of associated rewards

DQN: Experience Replay V

Initialise replay memory D to capacity N

Initialise action-value function Q with random weights

for episode = 1..M do

Initialise state s_t

for t = 1..T do

 Action a_t = select a random action a_t

 Execute action a_t and observe reward r and state s_{t+1}

 Store transition (s_t, a_t, r, s_{t+1}) in D

 Set $\theta_t = \theta$

 Sample random minibatch of transitions (s_t, a_t, r, s_{t+1}) from D

 Set $y_j = \begin{cases} r_j & \text{for terminal } s_{t+1} \\ r_j + \gamma \max_a Q(s_{t+1}, a'; \theta) & \text{for non-terminal } s_{t+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(s_t, a_t; \theta))^2$

 end for

end for

DQN: Target Network I

Whenever we run an update step on a Q-network's state we also update "near-by" states (monolithic architecture + generalisation ability of CNNs).

$$Q(s_t, a) \leftarrow Q(s_t, a) + \alpha[r_{t+1} + \gamma \max Q(s_{t+1}, a') - Q(s_t, a)] \quad (4)$$

Issue: Unstable training resulting from bootstrapping a continuous state space representation

- We may have an unstable learning process, because changes to $Q(s'_t, a)$ change $Q(s'_t, a')$ (with $s = s' \Rightarrow s'$) which changes $Q(s_t, a)$ on the next CNN and update in turn $Q(s'_t, a')$ forth.

- The effect may result in a run-away bias from bootstrapping and subsequently dominating the system numerically, causing the estimated Q values to diverge.

- In calculating the TD-error, the target function is changing too frequently when using convolutional neural networks.

Solution: Slow things down (as resonance damper)

Using a separate target network Q' , updated every few time (e.g. every 1000 steps) with a copy of the latest learned weight parameters, controls this stability (relaxation time).

- Initialising two Q networks: a main Q-network (Q), and a target network (Q')

When calculating the TD-error, use Q' , not Q

- Frequently set $Q = Q'$

- This gives the highly fluctuating Q-time to settle (we can measure this so called "relaxation time" empirically and set it accordingly) before updating Q'

Closing Remarks: Variation in rewards makes the training unstable → Clip positive reward to 1, and negative reward to -1.

DQN: Clipping rewards I

Clipping rewards:

- Each ATARI game has different score scales. For example, in Pong, players get 1 point when winning the play.

Otherwise, players get -1 point. However, in Space Invaders, players get 10-30 points when defeating invaders. This difference would make training unstable.

- Clipping Rewards: clips the rewards, with all positive rewards are set +1 and all negative rewards are set -1.

DQN: Skipping frames I

Skipping frames:

- Computer can simulate games (e.g. the ATARI simulator does this at 60 Hz) much faster than a human player would be able to react to the game, and this implies that the video game design and the required actions can be slower.
- Frame skipping uses only every 4 video game frames (i.e. 15 Hz...) (and frame stacking uses the past 4 frames as inputs), thus reducing computational cost and accelerating training times. It also makes the game run at a speed comparable to human reaction time.

Doubble Q network

Normal netowrk will produce a maximisation bias

Max at the Casino II

- So when update the value for $Q(Start, Left)$ we perform a max over actions in successor state Casino (i.e. from traces yielding r in Casino, A , Hotel) = -1, r (Casino, B , Hotel) = +1, r (Casino, A , Hotel) = -1) and immediately note down the +1 reward from choosing B .
- In the limit of large number of trials all actions will yield an average reward of -0.02, but that is only asymptotically reached.
- This results in a maximisation bias, namely that when taking a max over all actions with (very) finite data we may always overestimate values. This illustrates a general relationship (that we do not prove here), but we can see empirically.

Use the main network instead of the target network to do action selection (compared with normal target network)

Double Q Learning (van Hasselt et al., 2017, AAAI) I

Situation: In regular Deep-Q learning we use the target network for two tasks:

- Identify action with highest Q value,

- Compute Q value of this action

Complication: The maximum Q value may be overestimated (variance-bias problem), because an unusually high value from the main network Q does not mean that there is an unusually high value from the target network Q' . The problem with Deep-Q Learning is that the same samples (i.e. the Q network) are being used to decide which action is the best (highest expected reward), and the same samples are also being used to estimate that specific action-value.

Solution: In Double Q-Learning (van Hasselt et al., 2017, AAAI) we separate the two estimates: We use a target network Q' for 1. but use the regular Q for 2. (or the converse). Thus, we make the selection of the action with highest Q value, and selection of the Q value used in Bellman updates become independent from each other (different networks) so we become less susceptible to variance in each estimate. This reduces the frequency by which the maximum Q value may be overestimated, as it is less likely that both the networks are overestimating the same action.

REINFORCE was the crucial first step at popularising Policy Gradients (PGs)

REINFORCE algorithms suffer from high variance in the sampled trajectories, thus stabilising model parameters is difficult.

Direct Policy Gradients: Derivation idea 2/2

$$\theta^* = \arg \max_{\theta} J(\theta)$$

$$J(\theta) = E_{\pi \sim \pi(\cdot | s, \theta)}[r_t^\pi]$$

$$\log \text{of both sides}$$

$$\frac{1}{\theta} \log \pi(a_t | s_t, \theta) = \pi(a_t) \prod_{i=1}^T \pi(a_i | s_i, \theta)$$

$$\nabla_\theta J(\theta) = E_{\pi \sim \pi(\cdot | s, \theta)}[\nabla_\theta \log \pi(a_t | s_t, \theta)]$$

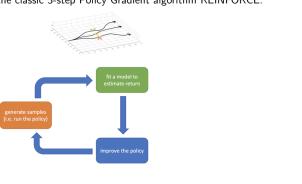
$$\nabla_\theta J(\theta) = E_{\pi \sim \pi(\cdot | s, \theta)} \left[\sum_{i=1}^T \nabla_\theta \log \pi(a_i | s_i, \theta) \right]$$

We arrived at the final equation, a version of the policy gradient theorem, by expanding the state sequence from a fixed starting point s_1 . Similarly we can obtain results for the stationary distribution of starting states (not discussed here).

REINFORCE algorithms

REINFORCE-ing good trajectories over bad ones

A direct implementation of the Policy Gradient Theorem result is the classic 3-step Policy Gradient algorithm REINFORCE.



• Before we had to learn a value function through function approximation and then derive a corresponding policy

- Often learning a value function can be intractable (more unstable / was at the time, intractable for large state spaces).

• REINFORCE provided a way to directly optimize policies to get around this problem.

REINFORCE algorithm:

$$1. \ sample \{ \tau \} \rightarrow \pi(a_t | s_t) \ (\text{run the policy})$$

$$2. \ \nabla_\theta J(\theta) \approx \sum_i (\sum_{t=1}^T \nabla_\theta \log \pi(a_t | s_t)) (\sum_{t=1}^T r(s_t, a_t))$$

$$3. \ \theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

• The method suffers from high variance in the sampled trajectories, thus stabilising model parameters is difficult.

What was the impact of REINFORCE? I

- Before: we had to learn a value function through function approximation and then derive a corresponding policy

- But, often learning a value function can be intractable (more unstable, computationally demanding for large state spaces, observability).

• REINFORCE provided a way to directly optimize policies to get around this problem. It was the direct result of rendering the PG theorem useful after its derivation

REINFORCE was the crucial first step at popularising Policy Gradients (PGs)

REINFORCE

- In the limit of large amounts of data, the model will converge to the optimal parameters

- But, the REINFORCE method suffers from high variance in the sampled trajectories, thus stabilising model parameters θ can be difficult.

- Any erratic trajectory can cause a sub-optimal shift in the policy distribution if by chance generates higher rewards or vice versa.

- Many subsequent algorithms were proposed to reduce the variance while keeping the bias unchanged by being smarter correlating rewards with trajectories. These may involve subtracting a baseline from the reward term (so as to keep values smaller) and by using a so called advantage term (Schulman et al., 2016).

Desired qualities for Reinforcement Learning

In robotics and many other real-world control systems it makes sense to have continuous policies. A simple way is to introduce Gaussian policies, i.e. $\pi(a_t | s_t) \sim \mathcal{N}(\mu(s_t), \Sigma)$ as a Gaussian distributions

$$\nabla_\theta \mu(\theta) = \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right)$$

$$\nabla_\theta \log \pi(a_t | s_t) = -\frac{1}{2} \sum_{i=1}^N (f(s_t) - \mu(s_t))^T \frac{\partial f}{\partial \theta}$$

We can interpret the term $f(s_t, a_t) - \mu(s_t)$ as a curious form of supervised learning – where we fit the network output computed on the states (s_t) to the actions a_t .

Policy gradients trial-and-error (like MC)

It is curious because pure maximum-likelihood-type supervised learning (fitting states and actions to each other) is actually reward-weighted fitting in Policy-Gradient-based reinforcement learning.

$$\nabla_\theta J(\theta) = \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi(a_t | s_t) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right)$$

$$\nabla_\theta J(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log \pi(a_t | s_t) \nabla_\theta \log \pi(a_t | s_t)$$

maximum likelihood: $\nabla_\theta \mu(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log \pi(a_t | s_t)$

These averages can be reinterpreted as weighted average over policies

- We increase the weight of trajectories that are good in reward

- We decrease the weight of trajectories that are bad in rewards

The action in the future cannot affect the actions in the past.

Ways to reduce variance: Causality I

- Apply Causality to reduce the amount of total data (and thus variability):

• Strictly Causality: $\nabla_\theta J(\theta) = \sum_{i=1}^N (\sum_{t=1}^T \nabla_\theta \log \pi(a_t | s_t)) (\sum_{t=1}^T r(s_t, a_t))$

but, the policy at time $t' < t$ cannot affect the reward in the past (time t), so we can restrict the terms entering the sums

$$\nabla_\theta J(\theta) = \sum_{i=1}^N (\sum_{t=1}^T \nabla_\theta \log \pi(a_t | s_t)) (\sum_{t=t'}^T r(s_t, a_t))$$

maximum likelihood: $\nabla_\theta \mu(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log \pi(a_t | s_t)$

• The parameter β controls the degree of blending (forgetting) between the current parameter estimate $\hat{\theta}$ and its update θ' at time t to determine the parameter vector at the next $t+1$.

Just like the Policy Gradient theorem can be thought of a policy-centered equivalent of the Bellman Theorem, DDPG can consider the continuous action version of the discrete action DQN. Why?

Desired qualities for Reinforcement Learning

• Stable learning & little sensitivity to hyperparameters: we want to avoid parameter tuning which is (computationally) expensive and more an art than a science. Ideally, RL algorithms work out of the box.

• Sample Efficiency: Good sample complexity is the prerequisite for efficient (and in some cases effective) skill acquisition

• Most algorithms we covered so far are relatively poor at sample efficiency (although we have seen introducing features improved their stability of learning). This is ok for simulation settings or "toy"-problems such as video games.

• Many real-world applications, e.g. robots require us to put this second aspect, good exploration at the centre of our attention, as interactive learning is costly and does not scale well.

An operational definition of the sample complexity of exploration for a reinforcement learning algorithm is the number of time steps in which the algorithm does not select near-optimal actions.

SAC – Objective Function I

In the Soft Actor-Critic we use the entropy concept to boost its sample complexity. The primary goal of SAC is to maximize a modified expected cumulative reward, which includes an entropy term to encourage exploration:

$$J(\pi) = E_{(s_t, a_t) \sim p_{\pi}} \left[\sum_t \gamma^t r(s_t, a_t) + \lambda H(\pi(\cdot | s_t)) \right] \quad (7)$$

Here,

- $J(\pi)$ is the cost-to-go objective function.

- $r(s_t, a_t)$ is the immediate reward function.

- $p_{\pi} = p_{\pi}(s, a)$ is the state-action distribution under policy π .

- γ is the discount factor.

- λ is a temperature parameter controlling the importance of the entropy term H .

- $\pi(\cdot | s_t)$ is the policy distribution given state s_t .

Why combine deep RL methods as Actor-Critics?

• Policy Gradient agents might take a lot of actions over the course of an episode, it is hard to assign credit to many actions, which means that these updates have a high variance

• Therefore, it may take a lot of updates for your policy to converge

• Moreover, PG-method only works in episodic regimes (cf. Monte Carlo tabular RL). If your agent is not acting in an episodic environment, it will never get an update (cf. "You cannot backup from death")

Directly calculating the policy gradient

- Generate the trace

- Update the policy value

- calculate the TD error

• update the Q value error

Naive, Deep, Q actor-critic (Weng, 2019) I

A simple Q-driven policy-gradient actor-critic.

Algorithm 1: Q-driven Actor-Critic

Initialise parameters θ , π , and learning rates $\alpha_\theta, \alpha_\pi$: sample $a \sim \pi(a | s, \theta)$.

Then sample $a' \sim \pi(a' | s', \theta)$ and state $s' \sim \pi(s' | s, a, \theta)$.

Update the next action $a'' \sim \pi(a'' | s'', \theta)$ and state $s'' \sim \pi(s'' | s', a', \theta)$.

Then update the next action $a''' \sim \pi(a''' | s''', \theta)$ and state $s''' \sim \pi(s''' | s'', a'', \theta)$.

... and use π to sample the next action of Q function:

$a'' \sim \pi(a'' | s'', a', \theta)$

More for $a \sim \pi(a | s, \theta)$ and $s' \sim \pi(s' | s, a, \theta)$.

end for

• update the Q value error

• update the policy gradient

• update the policy

• update the Q value

• update the policy gradient

• update the policy

• update the Q value