

Deep Q Network (DQN)

Announcements

We will release HW2 tonight (Q-learning, TD,
and simulation lemma)

We will release the first reading quiz today

Recap: Bellman operator

Value iteration

$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

R, \underline{P}

$$Q^{t+1} \leftarrow \underset{Q}{\operatorname{argmin}} \sum_{s \in \mathcal{S}, a \in \mathcal{A}} (Q(s, a) - (\mathcal{T} Q^t)(s, a))^2$$

$= 0$

Recap: Q-learning

Tabular Q Learning: maintain a table \hat{Q} of size $S \times A$

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \eta \left(r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) \right)$$

(s, a, r, s')

Data collection via ϵ -greedy:

Recap: Q-learning

Tabular Q Learning: maintain a table \hat{Q} of size $S \times A$

$$\hat{Q}(s, a) \leftarrow \hat{Q}(s, a) + \eta \left(r + \gamma \max_{a'} \hat{Q}(s', a') - \hat{Q}(s, a) \right)$$

Data collection via ϵ -greedy:

W/ prob ϵ , select action uniform randomly

W/ prob $1 - \epsilon$, select greedy action $\arg \max_a \hat{Q}(s, a)$

Today

Consider large-scale MDPs,

how to estimate $Q^*(s, a)$ using function approximation (e.g., neural network)

**Deep Q-network (DQN) is the earliest example
of showing Deep Learning + RL is powerful**

doi:10.1038/nature14236

Human-level control through deep reinforcement learning

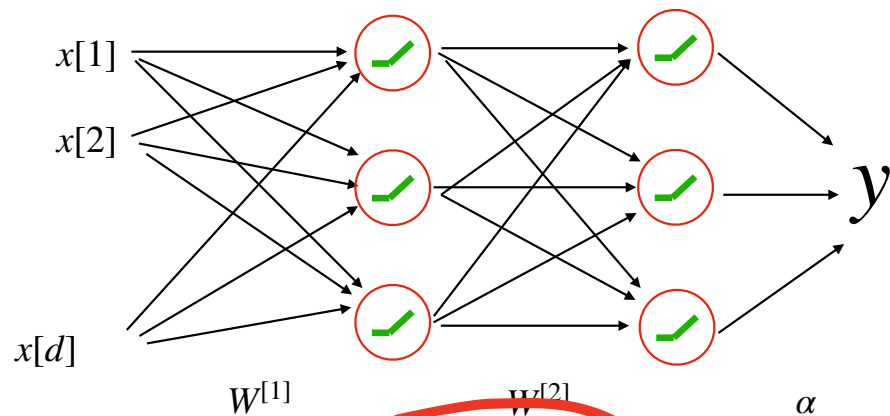
Volodymyr Mnih^{1*}, Koray Kavukcuoglu^{1*}, David Silver^{1*}, Andrei A. Rusu¹, Joel Veness¹, Marc G. Bellemare¹, Alex Graves¹, Martin Riedmiller¹, Andreas K. Fidjeland¹, Georg Ostrovski¹, Stig Petersen¹, Charles Beattie¹, Amir Sadik¹, Ioannis Antonoglou¹, Helen King¹, Dharshan Kumaran¹, Daan Wierstra¹, Shane Legg¹ & Demis Hassabis¹

Outline:

1. Q Learning w/ function approximation
2. Replay buffer, batch optimization and target network

Q-Learning w/ function approximation

We will model Q^* using a function approximator



$$Q_{\theta}(s, a) : S \times A \rightarrow \mathbb{R}$$

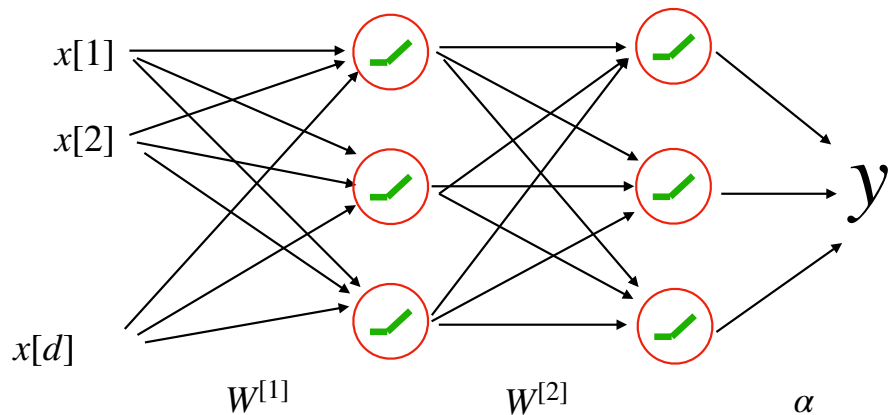
$$x = [s^{\top}, a]^{\top}$$

$$y = \alpha^{\top} \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$

$$\theta = [w^{[1]}, w^{[2]}, \alpha]$$

Q-Learning w/ function approximation

We will model Q^* using a function approximator



$$Q_{\theta}(s, a) : S \times A \rightarrow \mathbb{R}$$

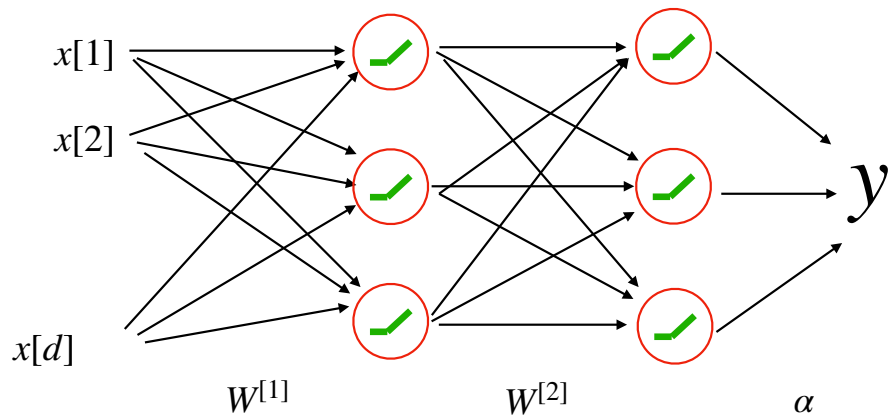
Assumption: differentiable $\nabla_{\theta} Q_{\theta}(s, a)$

$$x = [s^{\top}, a]^{\top}$$

$$y = \alpha^{\top} \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$

Q-Learning w/ function approximation

We will model Q^* using a function approximator



$$x = [s^\top, a]^\top$$

$$y = \alpha^\top \text{ReLU} \left(W^{[2]} \text{ReLU} \left(W^{[1]} x \right) \right) + b$$

Raw pixel-image

$$Q_\theta(s, a) : S \times A \rightarrow \mathbb{R}$$


Assumption: differentiable $\nabla_\theta Q_\theta(s, a)$

(The DQN paper uses ConvNet as s is an image frame of the game)

Attempt 1: Q Learning w/ function approximation

Initialize θ^0 . Set initial state $s \in \mathcal{S}$

For $t = 0$ to T



Attempt 1: Q Learning w/ function approximation

Initialize θ^0 . Set initial state $s \in \mathcal{S}$

For $t = 0$ to T

Take action a based on ϵ -greedy of Q_{θ^t} , get reward r and next state $s' \sim P(\cdot | s, a)$

env.step(a)

Attempt 1: Q Learning w/ function approximation

Initialize θ^0 . Set initial state $s \in \mathcal{S}$

For $t = 0$ to T

Take action a based on ϵ -greedy of Q_{θ_t} , get reward r and next state $s' \sim P(\cdot | s, a)$

Form Q-target $r + \gamma \max_{a'} Q_{\theta_t}(s', a')$

$Q_{\theta_t}(s, a)$ \rightarrow $r + \gamma \max_{a'} Q_{\theta_{t+1}}(s', a')$

Attempt 1: Q Learning w/ function approximation

Initialize θ^0 . Set initial state $s \in \mathcal{S}$

For $t = 0$ to T

Take action a based on ϵ -greedy of Q_{θ^t} , get reward r and next state $s' \sim P(\cdot | s, a)$

Form Q-target $r + \gamma \max_{a'} Q_{\theta^t}(s', a')$

Update to θ^{t+1} : \leftarrow SGD on BE loss

Attempt 1: Q Learning w/ function approximation

Initialize θ^0 . Set initial state $s \in \mathcal{S}$

For $t = 0$ to T

Take action a based on ϵ -greedy of Q_{θ^t} , get reward r and next state $s' \sim P(\cdot | s, a)$

Form Q-target $r + \gamma \max_{a'} Q_{\theta^t}(s', a')$

Update to θ^{t+1} :

Set $s \leftarrow s'$

Q Learning w/ function approximation

Update parameters using SGD on the Bellman error loss

Q Learning w/ function approximation

Update parameters using SGD on the Bellman error loss

$T \theta_\theta$

$\neq 0, \text{Hs.}$

$$\ell_{be}(\theta) := \underbrace{(Q_\theta(s, a) - y)^2}_{\text{predictor}}, \text{ where } y = r(s, a) + \underbrace{\gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q_\theta(s', a')}_{\text{Target}}$$

$$\nabla_{\theta} \ell_{BE}(\theta) = 2 \cdot (Q_\theta(s, a) - y) \cdot \nabla_{\theta} Q_\theta(s, a)$$

$$\tilde{\nabla}_{\theta} = 2 \cdot (Q_\theta(s, a) - (\underbrace{r + \gamma \max_{a'} Q_\theta(s', a')}_{\text{Target}})) \cdot \nabla_{\theta} Q_\theta(s, a)$$

$$\theta \leftarrow \theta - \eta \cdot \tilde{\nabla}_{\theta}$$

Issues of this simple approach

1. Inefficient — it throws away all historical data
(your network could forget old experiences, i.e., catastrophic forgetting)

Issues of this simple approach

1. Inefficient — it throws away all historical data
(your network could forget old experiences, i.e., catastrophic forgetting)

2. instability — Training is quite unstable (we saw it from the past Cartpole Demo)

Outline:

1. Q Learning w/ function approximation
2. Replay buffer, batch optimization and target network

Q-Learning w/ function approximation

(s, a, r, s')

First improvement: Replay buffer

$$\mathcal{D}_{rb} = \begin{bmatrix} \dots \\ (s, a, r, s') \\ \dots \\ \dots \end{bmatrix}$$

A dataset that contains all historical
State-action-reward-next state tuples

Q-Learning w/ function approximation

With replay buffer, we can use **mini-batch SGD** to update Bellman error loss

Q-Learning w/ function approximation

With replay buffer, we can use **mini-batch SGD** to update Bellman error loss

Given Q_{θ^t} and replay buffer \mathcal{D}_{rb} , randomly sample a mini-batch \mathcal{B} from \mathcal{D}_{rb}

$$\frac{1}{|\mathcal{B}|} \sum_{s, a, s'} \left(Q_{\theta}(s, a) - r - \max_{a'} Q_{\theta}(s, a') \right)^2$$

$\mathcal{B} = \{s, a, r, s'\}$
simple (s, a)

Q-Learning w/ function approximation

With replay buffer, we can use **mini-batch SGD** to update Bellman error loss

Given Q_{θ^t} and replay buffer \mathcal{D}_{rb} , randomly sample a mini-batch \mathcal{B} from \mathcal{D}_{rb}

$$\theta^{t+1} = \theta^t - \eta \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} \left(Q_{\theta^t}(s, a) - r - \max_{a'} Q_{\theta^t}(s', a') \right) \nabla_{\theta} Q_{\theta^t}(s, a)$$

mini-batch SGD

Q-Learning w/ function approximation

Second improvement: Target network (making Q learning more stable)

Q-Learning w/ function approximation

Second improvement: Target network (making Q learning more stable)

Recall that Q learning can be understood as running SGD on an **evolving** loss function

Q-Learning w/ function approximation

Second improvement: Target network (making Q learning more stable)

Recall that Q learning can be understood as running SGD on an **evolving** loss function

$$\ell_{be}(\theta) := \underbrace{\left(Q_{\theta}(s, a) - \underbrace{r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q_{\theta}(s', a')}_{\text{regression target}} \right)^2}_{\Delta}$$

Q-Learning w/ function approximation

Second improvement: Target network (making Q learning more stable)

Recall that Q learning can be understood as running SGD on an **evolving** loss function

$$\ell_{be}(\theta) := (Q_{\theta}(s, a) - y)^2, \text{ where } y = r(s, a) + \underbrace{\gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q_{\theta}(s', a')}_{\text{regression target}}$$

Source of instability: target changes immediately whenever we update θ

Q-Learning w/ function approximation

Second improvement: Target network (making Q learning more stable)

Introducing target network $Q_{\tilde{\theta}}$ to **slow down** the evolution of the BE loss

(e.g., set $\tilde{\theta}$ as a copy of an older version of θ)

Q-Learning w/ function approximation

Second improvement: Target network (making Q learning more stable)

Introducing target network $Q_{\tilde{\theta}}$ to **slow down** the evolution of the BE loss

(e.g., set $\tilde{\theta}$ as a copy of an older version of θ)

$$\ell_{be}(\theta) := (Q_{\theta}(s, a) - y)^2, \text{ where } y = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q_{\tilde{\theta}}(s', a')$$

Q-Learning w/ function approximation

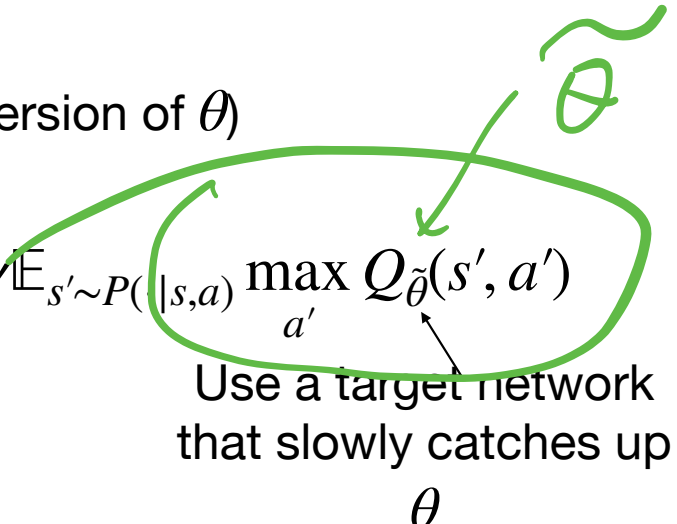
Second improvement: Target network (making Q learning more stable)

Introducing target network $Q_{\tilde{\theta}}$ to **slow down** the evolution of the BE loss

(e.g., set $\tilde{\theta}$ as a copy of an older version of θ)

$$\ell_{be}(\theta) := (Q_{\theta}(s, a) - y)^2, \text{ where } y = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a'} Q_{\tilde{\theta}}(s', a')$$

Use a target network
that slowly catches up
 θ



Attempt 2: Deep Q network (DQN)

Initialize θ and replay buffer \mathcal{D}_{rb} . Set $\tilde{\theta} = \theta$, Set initial state $s \in \mathcal{S}$

Q_θ \tilde{Q}_θ

While true:

Take action a based on ϵ -greedy of Q_θ , get reward r and next state $s' \sim P(\cdot | s, a)$

Attempt 2: Deep Q network (DQN)

Initialize θ and replay buffer \mathcal{D}_{rb} . Set $\tilde{\theta} = \theta$, Set initial state $s \in \mathcal{S}$

While true:

Take action a based on ϵ -greedy of Q_θ , get reward r and next state $s' \sim P(\cdot | s, a)$

Add (s, a, r, s') to \mathcal{D}_{rb}

Attempt 2: Deep Q network (DQN)

Initialize θ and replay buffer \mathcal{D}_{rb} . Set $\tilde{\theta} = \theta$, Set initial state $s \in \mathcal{S}$

While true:

Take action a based on ϵ -greedy of Q_{θ} , get reward r and next state $s' \sim P(\cdot | s, a)$

Add (s, a, r, s') to \mathcal{D}_{rb}

Sample mini-batch \mathcal{B} from \mathcal{D}_{rb}

Attempt 2: Deep Q network (DQN)

Initialize θ and replay buffer \mathcal{D}_{rb} . Set $\tilde{\theta} = \theta$, Set initial state $s \in \mathcal{S}$

While true:

Take action a based on ϵ -greedy of Q_θ , get reward r and next state $s' \sim P(\cdot | s, a)$

Add (s, a, r, s') to \mathcal{D}_{rb}

Sample mini-batch \mathcal{B} from \mathcal{D}_{rb}

Update parameters:

$$\theta \leftarrow \theta - \eta \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} \left(Q_\theta(s, a) - r - \max_{a'} Q_{\tilde{\theta}}(s', a') \right) \nabla_\theta Q_\theta(s, a)$$

"SG" on $\underline{\underline{\mathcal{B}}}$

$\tilde{\theta}$ Target Network

Attempt 2: Deep Q network (DQN)

Initialize θ and replay buffer \mathcal{D}_{rb} . Set $\tilde{\theta} = \theta$, Set initial state $s \in \mathcal{S}$

While true:

Take action a based on ϵ -greedy of Q_θ , get reward r and next state $s' \sim P(\cdot | s, a)$

Add (s, a, r, s') to \mathcal{D}_{rb}

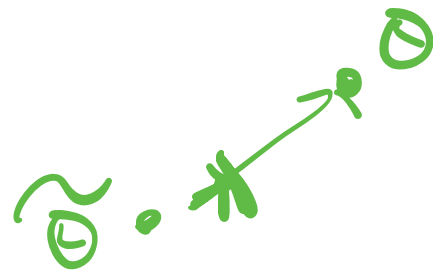
Sample mini-batch \mathcal{B} from \mathcal{D}_{rb}

Update parameters:

$$\theta \leftarrow \theta - \eta \frac{1}{|\mathcal{B}|} \sum_{(s,a,r,s') \in \mathcal{B}} \left(Q_\theta(s, a) - r - \max_{a'} Q_{\tilde{\theta}}(s', a') \right) \nabla_\theta Q_\theta(s, a)$$

Every C step, set $\tilde{\theta} = \theta$

Option 2: $\tilde{\theta} = (1-d) \cdot \tilde{\theta} + d \theta$



When C is large...

DQN is performing SGD for standard regression between two target network updates..

$$\min_{\theta} \sum_{s,a,r,s' \in \mathcal{D}_{rb}} \left(Q_{\theta}(s, a) - \left(r + \max_{a'} Q_{\tilde{\theta}}(s', a') \right) \right)^2 \approx 0$$

When C is large...

DQN is performing SGD for standard regression between two target network updates..

C-steps of SGD

$$\hat{\theta} \leftarrow \min_{\theta} \sum_{s,a,r,s' \in \mathcal{D}_{rb}} \left(Q_{\theta}(s,a) - \left(r + \max_{a'} Q_{\tilde{\theta}}(s',a') \right) \right)^2 \underbrace{E[y|s,a]}$$

Q: What is the Bayes optimal of this regression problem?

$$E[y|s,a] = r(s,a) + \gamma E_{s'|p(\cdot|s,a)} \max_{a'} Q_{\tilde{\theta}}(s',a')$$

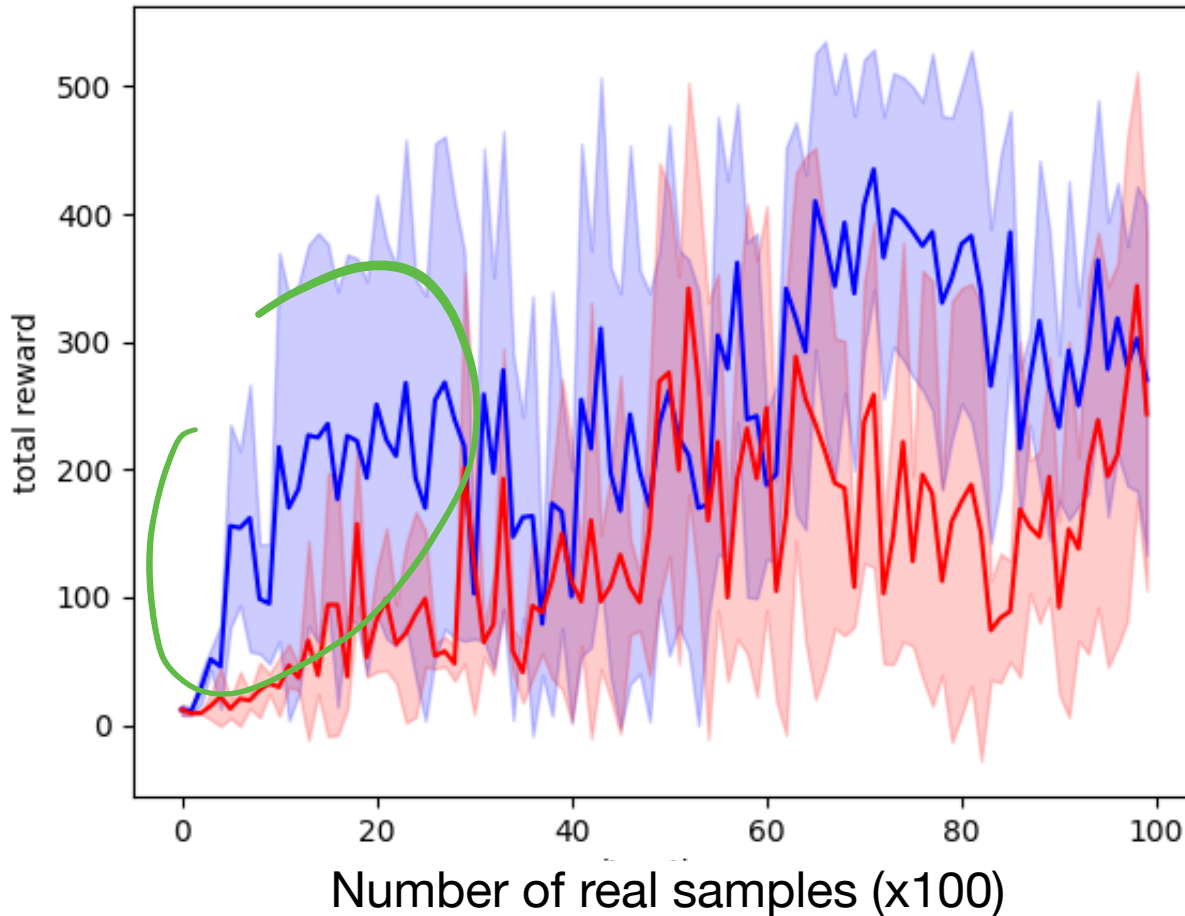
$$Q_{\hat{\theta}} \approx T Q_{\tilde{\theta}}$$

$$\Rightarrow \tilde{\theta} = \hat{\theta}, \text{ repeat}$$

$$= (T Q_{\hat{\theta}})(s,a)$$

DQN vs Naive Q-learning

DQN (blue) vs Q-learning (red)



$c=5$

BLUE: DQN

RED: Q-learning

Summary

1. Using function approximation to handle large state space
2. Making Q-learning closer to the supervised learning (i.e., regression) framework:
 - Replay buffer + mini-batch SGD
 - Target network to simulate a standard regression setting