RL + DL

Rob Platt Northeastern University

Some slides from David Silver, Deep RL Tutorial

Deep RL in an exciting recent development



Deep RL in an exciting recent development



Remember the RL scenario...



Tabular RL

until S is terminal



Tabular RL



Deep RL



Deep RL



Where does "state" come from?



Earlier, we dodged this question: "it's part of the MDP problem statement"

But, that's a cop out. How do we get state?

Typically can't use "raw" sensor data as state w/ a tabular Q-function – it's too big (e.g. pacman has something like 2^(num pellets) + ... states)

Linear function approximation?



In the RL lecture, we suggested using Linear Function approximation:

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \dots + w_n f_n(s,a)$$

Q-function encoded by the weights

Linear function approximation?



In the RL lecture, we suggested using Linear Function approximation:

$$Q(s,a) = w_1 f_1(s,a) + w_2 f_2(s,a) + \dots + w_n f_n(s,a)$$

Q-function encoded by the weights

But, this STILL requires a human designer to specify the features!

Linear function approximation?





Deep Q Learning



Convolutional Agent

Deep Q Learning



Instead of state, we have an image

 in practice, it could be a history of the k most recent images stacked as a single k-channel image

Hopefully this new image representation is Markov...

– in some domains, it might not be!

Network structure





Network structure



Here's the standard Q-learning update equation:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

Here's the standard Q-learning update equation:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

Initialize $Q(s, a), \forall s \in S, a \in \mathcal{A}(s)$, arbitrarily, and $Q(terminal-state, \cdot) = 0$ Repeat (for each episode):

Initialize S

Repeat (for each step of episode):

Choose A from S using policy derived from Q (e.g., ϵ -greedy) Take action A, observe R, S'

$$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_{a} Q(S', a) - Q(S, A)]$$

$$S \leftarrow S'$$

until S is terminal

Here's the standard Q-learning update equation:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

Rewriting:

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a')\right]$$

Here's the standard Q-learning update equation:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

Rewriting:

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a')\right]$$

let's call this the "target"

This equation adjusts Q(s,a) in the direction of the target

Here's the standard Q-learning update equation:

$$Q(s,a) \leftarrow Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a') - Q(s,a) \right]$$

Rewriting:

$$Q(s,a) \leftarrow (1-\alpha)Q(s,a) + \alpha \left[r + \gamma \max_{a'} Q(s',a')\right]$$

let's call this the "target"

This equation adjusts Q(s,a) in the direction of the target

We're going to accomplish this same thing in a different way using neural networks...

Remember how we train a neural network...

Given a dataset:
$$D = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$$

Define loss function: $L(x^i, y^i; w, b) = \frac{1}{2}(y^i - f(w^T x^i + b))^2$

Adjust *w*, *b* so as to minimize:
$$\sum_{(x^i, y^i) \in D} L(x^i, y^i; w, b)$$

Do gradient descent on dataset:

1. repeat 2. $w \leftarrow w - \alpha \frac{1}{n} \sum_{(x^i, y^i) \in D} \nabla_w L(x^i, y^i; w, b)$ 3. $b \leftarrow b - \alpha \frac{1}{n} \sum_{(x^i, y^i) \in D} \nabla_b L(x^i, y^i; w, b)$

4. until converged

Use this loss function:

$$L(s, a, s'; w) = \frac{1}{2} \left(r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right)^2$$

Use this loss function:

$$L(s, a, s'; w) = \frac{1}{2} \left(r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right)^2$$
Notice that Q is now

Notice that Q is now parameterized by the weights, *w*



Use this loss function: $\operatorname{target} L(s, a, s'; w) = \frac{1}{2} \left(r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right)^2$





Now we can do gradient descent!

Deep Q Learning, no replay buffers

Initialize Q(s,a;w) with random weights Repeat (for each episode):

Initialize s

Repeat (for each step of the episode):

Choose *a* from *s* using policy derived from Q (e.g. e-greedy) Take action *a*, observe *r*, *s*'

$$w \leftarrow w - \alpha \nabla_w L(s, a, s'; w)$$
$$s \leftarrow s'$$

Until s is terminal

Where:

$$\nabla_w L(s, a, s'; w) \doteq -\left(r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a)\right) \nabla_w Q_w(s, a)$$

Deep Q Learning, no replay buffers



Example: 4x4 frozen lake env

Get to the goal (G) Don't fall in a hole (H)

FHFH FFFH HFFG (Left)

steps:	11582,	episodes:	770,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0873818397522
steps:	11609,	episodes:	771,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0872020721436
steps:	11652,	episodes:	772,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.138998985291
steps:	11672,	episodes:	773,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0649240016937
steps:	11689,	episodes:	774,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0546970367432
steps:	11697,	episodes:	775,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0260739326477
steps:	11731,	episodes:	776,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.110991954803
steps:	11773,	episodes:	777,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.135339975357
steps:	11798,	episodes:	778,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0810689926147
steps:	11818,	episodes:	779,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0643260478973
steps:	11870,	episodes:	780,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.169064044952
steps:	11906,	episodes:	781,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.117113113403
steps:	11992,	episodes:	782,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.279519796371
steps:	12064,	episodes:	783,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.234206199646
steps:	12090,	episodes:	784,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.0835938453674
steps:	12137,	episodes:	785,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.150979042053
steps:	12185,	episodes:	786,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.155304908752
steps:	12245,	episodes:	787,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.194122076035
steps:	12277,	episodes:	788,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.102608919144
steps:	12293,	episodes:	789,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.0520431995392

Demo!

Experience replay

Deep learning typically assumes independent, identically distributed (IID) training data

Experience replay

Deep learning typically assumes independent, identically distributed (IID) training data

```
But is this true in the deep RL scenario?
```

```
Initialize Q(s,a;w) with random weights Repeat (for each episode):
```

Initialize s

Repeat (for each step of the episode):

```
Choose a from s using policy derived from Q (e.g. e-greedy)
Take action a, observe r, s'
```

```
w \leftarrow w - \alpha \nabla_w L(s, a, s'; w)s \leftarrow s'
```

Until s is terminal

Experience replay

Deep learning typically assumes independent, identically distributed (IID) training data

But is this true in the deep RL scenario?

```
Initialize Q(s,a;w) with random weights
Repeat (for each episode):
Initialize s
Repeat (for each step of the episode):
Choose a from s using policy derived from Q (e.g. e-greedy)
```

Our solution: buffer experiences and then "replay" them during training

Deep Q Learning WITH replay buffers



Deep Q Learning WITH replay buffers

$$L(B;w) = \frac{1}{2} \sum_{(s,a,s',r)\in B} \left(r + \gamma \max_{a'} Q_w(s',a') - Q_w(s,a)\right)^2$$

- buffers like this are pretty common in DL...
$$s \leftarrow s$$

If mod(step,transpace) = 0: Train every
sample batch from D
w \leftarrow w - \alpha \nabla_w L(B;w)
One step grad
descent WRT buffer

Example: 4x4 frozen lake env

Get to the goal (G) Don't fall in a hole (H)

FHFH FFFH HFFG (Left)

steps:	11582,	episodes:	770,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0873818397522
steps:	11609,	episodes:	771,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0872020721436
steps:	11652,	episodes:	772,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.138998985291
steps:	11672,	episodes:	773,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0649240016937
steps:	11689,	episodes:	774,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0546970367432
steps:	11697,	episodes:	775,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0260739326477
steps:	11731,	episodes:	776,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.110991954803
steps:	11773,	episodes:	777,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.135339975357
steps:	11798,	episodes:	778,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0810689926147
steps:	11818,	episodes:	779,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.0643260478973
steps:	11870,	episodes:	780,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.169064044952
steps:	11906,	episodes:	781,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.117113113403
steps:	11992,	episodes:	782,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.279519796371
steps:	12064,	episodes:	783,	mean	100	episode	reward:	0.6,	%	time	spent	exploring:	2,	time	elapsed:	0.234206199646
steps:	12090,	episodes:	784,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.0835938453674
steps:	12137,	episodes:	785,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.150979042053
steps:	12185,	episodes:	786,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.155304908752
steps:	12245,	episodes:	787,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.194122076035
steps:	12277,	episodes:	788,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.102608919144
steps:	12293.	episodes:	789,	mean	100	episode	reward:	0.7,	%	time	spent	exploring:	2,	time	elapsed:	0.0520431995392

Demo!

Atari Learning Environment



Atari Learning Environment

- End-to-end learning of values Q(s, a) from pixels s
- Input state s is stack of raw pixels from last 4 frames
- Output is Q(s, a) for 18 joystick/button positions
- Reward is change in score for that step



Network architecture and hyperparameters fixed across all games

Atari Learning Environment



- **Double DQN**: Remove upward bias caused by max $Q(s, a, \mathbf{w})$
 - Current Q-network w is used to select actions
 - Older Q-network w⁻ is used to evaluate actions

$$I = \left(r + \gamma Q(s', \operatorname{argmax}_{a'} Q(s', a', \mathbf{w}), \mathbf{w}^{-}) - Q(s, a, \mathbf{w})\right)^{2}$$

- **Double DQN**: Remove upward bias caused by max $Q(s, a, \mathbf{w})$
 - Current Q-network w is used to select actions
 - ► Older Q-network **w**⁻ is used to evaluate actions

$$I = \left(r + \gamma Q(s', \operatorname{argmax}_{a'} Q(s', a', \mathbf{w}), \mathbf{w}^{-}) - Q(s, a, \mathbf{w})\right)^{2}$$

- Prioritised replay: Weight experience according to surprise
 - Store experience in priority queue according to DQN error

$$\left| r + \gamma \max_{a'} Q(s', a', \mathbf{w}^{-}) - Q(s, a, \mathbf{w}) \right|$$

- **Double DQN**: Remove upward bias caused by max $Q(s, a, \mathbf{w})$
 - Current Q-network w is used to select actions
 - Older Q-network \mathbf{w}^- is used to evaluate actions

$$I = \left(r + \gamma Q(s', \operatorname{argmax}_{a'} Q(s', a', \mathbf{w}), \mathbf{w}^{-}) - Q(s, a, \mathbf{w})\right)^{2}$$

- Prioritised replay: Weight experience according to surprise
 - Store experience in priority queue according to DQN error

$$\left| \mathbf{r} + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}', \mathbf{w}^{-}) - Q(\mathbf{s}, \mathbf{a}, \mathbf{w}) \right|$$

- Duelling network: Split Q-network into two channels
 - Action-independent value function V(s, v)
 - Action-dependent advantage function A(s, a, w)

$$Q(s,a) = V(s,v) + A(s,a,w)$$

- **Double DQN**: Remove upward bias caused by max $Q(s, a, \mathbf{w})$
 - Current Q-network w is used to select actions
 - Older Q-network w⁻ is used to evaluate actions

$$I = \left(r + \gamma Q(s', \operatorname{argmax}_{a'} Q(s', a', \mathbf{w}), \mathbf{w}^{-}) - Q(s, a, \mathbf{w})\right)^{2}$$

- Prioritised replay: Weight experience according to surprise
 - Store experience in priority queue according to DQN error

$$\left| \mathbf{r} + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}', \mathbf{w}^{-}) - Q(\mathbf{s}, \mathbf{a}, \mathbf{w}) \right|$$

- Duelling network: Split Q-network into two channels
 - Action-independent value function V(s, v)
 - Action-dependent advantage function A(s, a, w)

$$Q(s,a) = V(s,v) + A(s,a,\mathbf{w})$$

Combined algorithm: 3x mean Atari score DQN w/ replay buffers