## **DEEP LEARNING**

# Introduction to Deep Reinforcement Learning

deeplearning.mit.edu

2019

## 最专业报告分享群:

#### •每日分享5+科技行业报告

- 同行业匹配,覆盖人工智能、大数据、机器人、 智慧医疗、智能家居、物联网等行业。
- 高质量用户,同频的人说同样的话

扫描右侧二维码, 或直接搜索关注公众号: 智东西(zhidxcom) 回复"报告群"加入



## Deep Reinforcement Learning (Deep RL) Deep Learning Deep RL



- What is it? Framework for learning to solve sequential decision making problems.
- How? Trial and error in a world that provides occasional rewards
- **Deep?** Deep RL = RL + Neural Networks

## **Types of Learning**

- Supervised Learning
- Semi-Supervised Learning
- Unsupervised Learning
- Reinforcement Learning



## It's all "supervised" by a loss function!

\*Someone has to say what's good and what's bad (see Socrates, Epictetus, Kant, Nietzsche, etc.)





Supervised learning is "teach by example":

Here's some examples, now learn patterns in these example.

Reinforcement learning is "teach by experience": Here's a world, now learn patterns by exploring it.



## Machine Learning: Supervised vs Reinforcement

Supervised learning is "teach by example": Here's some examples, now learn patterns in these example.

Reinforcement learning is "teach by experience": Here's a world, now learn patterns by exploring it.



### Success

## **Reinforcement Learning in Humans**

- Human appear to learn to walk through "very few examples" of trial and error. **How** is an open question...
- Possible answers:
  - Hardware: 230 million years of bipedal movement data.
  - Imitation Learning: Observation of other humans walking.
  - Algorithms: Better than backpropagation and stochastic gradient descent







### Open Question: What can be learned from data?









Lidar



Microphone

Camera (Visible, Infrared)

3





GPS





IMU

#### https://deeplearning.mit.edu 2019

Networking

(Wired, Wireless)

References: [132]

achusett

Institute of Technology



Institute of Technology







**Fechnology** 

#### Image Recognition: If it looks like a duck



#### Audio Recognition: Quacks like a duck



Activity Recognition: Swims like a duck





termat's equation: This equation has no solutions in integers for  $n \ge 3$ .

Final **breakthrough**, 358 years after its conjecture: "It was so indescribably beautiful; it was so simple and so elegant. I couldn't understand how I'd missed it and I just stared at it in disbelief for twenty minutes. Then during the day I walked around the department, and I'd keep coming back to my desk looking to see if it was still there. It was still there. I couldn't contain myself, I was so excited. It was the most important moment of my working life. Nothing I ever do again will mean as much."







Institute of Technology References: [133]





## **Reinforcement Learning Framework**

At each step, the agent:

- Executes action
- Observe new state
- Receive reward

#### **Open Questions:**

- What cannot be modeled in this way?
- What are the challenges of learning in this framework?





For the full list of references visit: [80] <u>https://hcai.mit.edu/references</u>

## **Environment and Actions**

- Fully Observable (Chess) vs Partially Observable (Poker)
- Single Agent (Atari) vs Multi Agent (DeepTraffic)
- Deterministic (Cart Pole) vs Stochastic (DeepTraffic)
- Static (Chess) vs Dynamic (DeepTraffic)
- Discrete (Chess) vs Continuous (Cart Pole)

**Note:** Real-world environment might not technically be stochastic or partially-observable but might as well be treated as such due to their complexity.



## The Challenge for RL in Real-World Applications

Reminder:

Supervised learning: teach by example

Reinforcement learning: teach by experience

### Open Challenges. Two Options:

- 1. Real world observation + one-shot trial & error
- 2. Realistic simulation + transfer learning



For the full list of references visit: <u>https://hcai.mit.edu/references</u>

## Major Components of an RL Agent

An RL agent may be directly or indirectly trying to learn a:

- **Policy:** agent's behavior function
- Value function: how good is each state and/or action
- Model: agent's representation of the environment

## Meaning of Life for RL Agent: Maximize Reward

- Future reward:  $R_t = r_t + r_{t+1} + r_{t+2} + \dots + r_n$
- Discounted future reward:

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{n-t} r_n$$

- A good strategy for an agent would be to always choose an action that maximizes the (discounted) future reward
- Why "discounted"?
  - Math trick to help analyze convergence
  - Uncertainty due to environment stochasticity, partial observability, or that life can end at any moment:

"If today were the last day of my life, would I want to do what I'm about to do today?" – Steve Jobs

## Robot in a Room



- We can learn the value of (action, state) pairs and act greed/non-greedy
- We can learn the policy directly while sampling from it

## **Optimal Policy for a Deterministic World**

## Reward: -0.04 for each step



actions: UP, DOWN, LEFT, RIGHT

When actions are deterministic:

UP

100%	move UP
0%	move LEFT
0%	move RIGHT

Policy: Shortest path.



## Reward: -0.04 for each step



actions: UP, DOWN, LEFT, RIGHT

When actions are stochastic:

UP

80%	move UP
10%	move LEFT
10%	move RIGHT

**Policy:** Shortest path. Avoid -UP around -1 square.



## Reward: -2 for each step



actions: UP, DOWN, LEFT, RIGHT

When actions are stochastic:

UP

80%	move UP
10%	move LEFT
10%	move RIGHT

Policy: Shortest path.



Reward: -0.1 for each step



Reward: -0.04 for each step



### More urgent

Less urgent



For the full list of references visit: <u>https://hcai.mit.edu/references</u>

https://deeplearning.mit.edu 2019

## Reward: +0.01 for each step



actions: UP, DOWN, LEFT, RIGHT

When actions are stochastic:

UP

80%	move UP
10%	move LEFT
10%	move RIGHT

Policy: Longest path.

### Lessons from Robot in Room

- Environment model has big impact on optimal policy
- Reward structure has big impact on optimal policy



Massachuset Institute of Technology

For the full list of references visit: <u>https://hcai.mit.edu/references</u>

## Reward structure may have Unintended Consequences

#### Human



#### AI (Deep RL Agent)



Player gets reward based on:

- 1. Finishing time
- 2. Finishing position
- 3. Picking up "turbos"

## AI Safety

#### Risk (and thus Human Life) Part of the Loss Function



Massachusetts Institute of Technology

For the full list of references visit: https://hcai.mit.edu/references



### **Cart-Pole Balancing**

- **Goal** Balance the pole on top of a moving cart
- **State** Pole angle, angular speed. Cart position, horizontal velocity.
- Actions horizontal force to the cart
- **Reward** 1 at each time step if the pole is upright

### Doom\*

- Goal: Eliminate all opponents
- State: Raw game pixels of the game
- Actions: Up, Down, Left, Right, Shoot, etc.
- Reward:
- Positive when eliminating an opponent, negative when the agent is eliminated

\* Added for important thought-provoking considerations of AI safety in the context of autonomous weapons systems (see AGI lectures on the topic).









#### Grasping Objects with Robotic Arm

- **Goal** Pick an object of different shapes
- State Raw pixels from camera
- Actions Move arm. Grasp.
- Reward Positive when pickup is successful







#### Human Life

- Goal Survival? Happiness?
- State Sight. Hearing. Taste. Smell. Touch.
- Actions Think. Move.
- Reward Homeostasis?



## Key Takeaways for Real-World Impact

- Deep Learning:
  - Fun part: Good algorithms that learn from data.
  - Hard part: Good questions, huge amounts of representative data.
- Deep Reinforcement Learning:
  - Fun part: Good algorithms that learn from data.
  - Hard part: Defining a useful state space, action space, and reward.
  - Hardest part: Getting meaningful data for the above formalization.





For the full updated list of references visit: https://selfdrivingcars.mit.edu/references

## **3 Types of Reinforcement Learning**



## **Model-based**

- Learn the model of the world, then plan using the model
- Update model often
- Re-plan often

### Value-based

- Learn the state or state-action value
- Act by choosing best action in state
- Exploration is a necessary add-on

## **Policy-based**

- Learn the stochastic policy function that maps state to action
- Act by sampling policy
- Exploration is baked in

## Taxonomy of RL Methods



Link: https://spinningup.openai.com


# Taxonomy of RL Methods





# **Q-Learning**

- State-action value function: Q<sup>π</sup>(s,a)
  - Expected return when starting in s, performing a, and following π



- Q-Learning: Use **any policy** to estimate Q that maximizes future reward:
  - Q directly approximates Q\* (Bellman optimality equation)
  - Independent of the policy being followed
  - Only requirement: keep updating each (s,a) pair





# **Exploration vs Exploitation**

- Deterministic/greedy policy won't explore all actions
  - Don't know anything about the environment at the beginning
  - Need to try all actions to find the optimal one
- ε-greedy policy
  - With probability 1-ε perform the optimal/greedy action, otherwise random action
  - Slowly move it towards greedy policy:  $\epsilon \rightarrow 0$







#### **Q-Learning: Value Iteration**



	A1	A2	A3	A4
S1	+1	+2	-1	0
S2	+2	0	+1	-2
S3	-1	+1	0	-2
S4	-2	0	+1	+1

```
initialize Q[num_states,num_actions] arbitrarily
observe initial state s
repeat
    select and carry out an action a
    observe reward r and new state s'
    Q[s,a] = Q[s,a] + \alpha(r + \gamma \max_{a'} Q[s',a'] - Q[s,a])
    s = s'
```

```
until terminated
```

### Q-Learning: Representation Matters

- In practice, Value Iteration is impractical
  - Very limited states/actions
  - Cannot generalize to unobserved states



- State: screen pixels
  - Image size: 84 × 84 (resized)
  - Consecutive **4** images
  - Grayscale with **256** gray levels

 $256^{84 \times 84 \times 4}$  rows in the Q-table!

 $= 10^{69,970} >> 10^{82}$  atoms in the universe



#### **Deep RL** = RL + Neural Networks



#### DQN: Deep Q-Learning

Use a neural network to approximate the Q-function:

 $Q(s,a;\boldsymbol{\theta}) \approx Q^*(s,a)$ 









# Deep Q-Network (DQN): Atari



Layer	Input	Filter size	Stride	Num filters	Activation	Output
conv1	84x84x4	8x8	4	32	ReLU	20x20x32
conv2	20x20x32	4x4	2	64	ReLU	9x9x64
conv3	9x9x64	3x3	1	64	ReLU	7x7x64
fc4	7x7x64			512	ReLU	512
fc5	512			18	Linear	18

Mnih et al. "Playing atari with deep reinforcement learning." 2013.



DQN and Double DQN

Loss function (squared error):

$$L = \mathbb{E}[(\mathbf{r} + \boldsymbol{\gamma} \max_{a'} \mathbf{Q}(s', a') - Q(s, a))^2]$$
  
target prediction

- DQN: same network for both Q
- Double DQN: separate network for each Q
  - Helps reduce bias introduced by the inaccuracies of Q network at the beginning of training



# **DQN** Tricks

- Experience Replay
  - Stores experiences (actions, state transitions, and rewards) and creates mini-batches from them for the training process
- Fixed Target Network
  - Error calculation includes the target function depends on network parameters and thus changes quickly. Updating it only every 1,000 steps increases stability of training process.

$$Q(s_t, a) \leftarrow Q(s_t, a) + lpha \left[ r_{t+1} + \gamma \max_p Q(s_{t+1}, p) - Q(s_t, a) 
ight]$$

Replay	0	0	×	×
Target	0	×	0	×
Breakout	316.8	240.7	10.2	3.2
River Raid	7446.6	4102.8	2867.7	1453.0
Seaquest	2894.4	822.6	1003.0	275.8
Space Invaders	1088.9	826.3	373.2	302.0

target Q function in the red rectangular is fixed



#### Atari Breakout



After **10 Minutes** of Training

After **120** Minutes of Training

After 240 Minutes of Training





#### **DQN** Results in Atari

Institute of Technology [83]

# Dueling DQN (DDQN)



Decompose Q(s,a)

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

- V(s): the value of being at that state
- **A(s,a)**: the **advantage** of taking action **a** in state **s** versus all other possible actions at that state
- Use two streams:
  - one that estimates the state value V(s)

[336]

- one that estimates the advantage for each action A(s,a)
- Useful for states where action choice does not affect Q(s,a)

# **Prioritized Experience Replay**



• Sample experiences based on impact not frequency of occurrence

[336]

# Taxonomy of RL Methods





# Policy Gradient (PG)

- DQN (off-policy): Approximate Q and infer optimal policy
- PG (on-policy): Directly optimize policy space



Policy Network

Good illustrative explanation: http://karpathy.github.io/2016/05/31/rl/

*"Deep Reinforcement Learning: Pong from Pixels"* 

# Policy Gradient – Training

Policy Gradients: Run a policy for a while. See what actions led to high rewards. Increase their probability.



• **REINFORCE:** Policy gradient that increases probability of good actions and decreases probability of bad action:

$$abla_ heta E[R_t] = E[
abla_ heta log P(a)R_t]$$



# **Policy Gradient**

- Pros vs DQN:
  - Messy World: If Q function is too complex to be learned, DQN may fail miserably, while PG will still learn a good policy.
  - Speed: Faster convergence
  - Stochastic Policies: Capable of learning stochastic policies DQN can't
  - **Continuous actions:** Much easier to model continuous action space
- Cons vs DQN:
  - Data: Sample inefficient (needs more data)
  - **Stability:** Less stable during training process.
  - Poor credit assignment to (state, action) pairs for delayed rewards

#### **Problem with REINFORCE:**

Calculating the reward at the end, means all the actions will be averaged as good because the total reward was high.



# Taxonomy of RL Methods





#### Advantage Actor-Critic (A2C)

- Combine DQN (value-based) and REINFORCE (policy-based)
- Two neural networks (Actor and Critic):
  - Actor is policy-based: Samples the action from a policy
  - Critic is value-based: Measures how good the chosen action is

Policy Update: 
$$\Delta \theta = \alpha * \nabla_{\theta} * (\log \pi(S_t, A_t, \theta)) * R(t)$$
  
New update:  $\Delta \theta = \alpha * \nabla_{\theta} * (\log \pi(S_t, A_t, \theta)) * Q(S_t, A_t)$ 

• Update at each time step - temporal difference (TD) learning



# Asynchronous Advantage Actor-Critic (A3C)



- Both use parallelism in training
- A2C syncs up for global parameter update and then start each iteration with the same policy

# Taxonomy of RL Methods





# Deep Deterministic Policy Gradient (DDPG)

- Actor-Critic framework for learning a deterministic policy
- Can be thought of as: DQN for continuous action spaces
- As with all DQN, following tricks are required:
  - Experience Replay
  - Target network
- Exploration: add noise to actions, reducing scale of the noise as training progresses



# Taxonomy of RL Methods





# **Policy Optimization**

- Progress beyond Vanilla Policy Gradient:
  - Natural Policy Gradient
  - TRPO
  - PPO
- Basic idea in on-policy optimization:

Avoid taking bad actions that collapse the training performance.





#### **Trust Region:** First pick step size, then direction







# Taxonomy of RL Methods





#### Game of Go



Game size	Board size N	3 <sup>N</sup>	Percent legal	legal game positions (A094777) <sup>[11]</sup>
1×1	1	3	33%	1
2×2	4	81	70%	57
3×3	9	19,683	<mark>64</mark> %	12,675
4×4	16	43,046,721	56%	24,318,165
5×5	25	8.47×10 <sup>11</sup>	49%	4.1×10 <sup>11</sup>
9×9	81	4.4×10 <sup>38</sup>	23.4%	1.039×10 <sup>38</sup>
13×13	169	4.3×10 <sup>80</sup>	8. <mark>6</mark> 6%	3.72497923×10 <sup>79</sup>
19×19	361	1.74×10 <sup>172</sup>	1.196%	2.08168199382×10 <sup>170</sup>

For the full updated list of references visit: [170] https://selfdrivingcars.mit.edu/references

#### AlphaGo (2016) Beat Top Human at Go





Massachusetts Institute of Technology

For the full updated list of references visit: https://selfdrivingcars.mit.edu/references [83] https://deeplearning.mit.edu 2019

### AlphaGo Zero (2017): Beats AlphaGo





#### AlphaGo Zero Approach

- Same as the best before: Monte Carlo Tree Search (MCTS)
  - Balance exploitation/exploration (going deep on promising positions or exploring new underplayed positions)
- Use a neural network as "intuition" for which positions to expand as part of MCTS (same as AlphaGo)



For the full updated list of references visit: <u>https://selfdrivingcars.mit.edu/references</u> [170]

# AlphaGo Zero Approach

- Same as the best before: Monte Carlo Tree Search (MCTS)
  - Balance exploitation/exploration (going deep on promising positions or exploring new underplayed positions)
- Use a neural network as "intuition" for which positions to expand as part of MCTS (same as AlphaGo)
- "Tricks"
  - Use MCTS intelligent look-ahead (instead of human games) to improve value estimates of play options
  - Multi-task learning: "two-headed" network that outputs (1) move probability and (2) probability of winning.
  - Updated architecture: use residual networks



#### AlphaZero vs Chess, Shogi, Go



#### AlphaZero vs Chess, Shogi, Go



#### AlphaZero vs Chess, Shogi, Go





For the full list of references visit: [340] https://hcai.mit.edu/references

#### To date, for most successful robots operating in the real world: Deep RL is not involved





#### To date, for most successful robots operating in the real world: Deep RL is not involved



Massachusetts Institute of Technology

For the full list of references visit: <u>https://hcai.mit.edu/references</u>
### But... that's slowly changing: Learning Control Dynamics





## But... that's slowly changing: Learning to Drive: Beyond Pure Imitation (Waymo)





# The Challenge for RL in Real-World Applications

Reminder:

Supervised learning: teach by example

Reinforcement learning: teach by experience

#### Open Challenges. Two Options:

- 1. Real world observation + one-shot trial & error
- 2. Realistic simulation + transfer learning



For the full list of references visit: <u>https://hcai.mit.edu/references</u>

### Thinking Outside the Box: Multiverse Theory and the Simulation Hypothesis

• Create an (infinite) set of simulation environments to learn in so that our reality becomes just another sample from the set.





For the full list of references visit: <u>https://hcai.mit.edu/references</u> [342]

# Next Steps in Deep RL

- Lectures: <u>https://deeplearning.mit.edu</u>
- Tutorials: <u>https://github.com/lexfridman/mit-deep-learning</u>
- Advice for Research (from <u>Spinning Up as a Deep RL Researcher</u> by Joshua Achiam)
  - Background
    - Fundamentals in probability, statistics, multivariate calculus.
    - Deep learning basics
    - Deep RL basics
    - TensorFlow (or PyTorch)
  - Learn by doing
    - Implement core deep RL algorithms (discussed today)
    - Look for tricks and details in papers that were key to get it to work
    - Iterate fast in simple environments (see Einstein quote on simplicity)
  - Research
    - Improve on an existing approach
    - Focus on an unsolved task / benchmark
    - Create a new task / problem that hasn't been addressed with RL