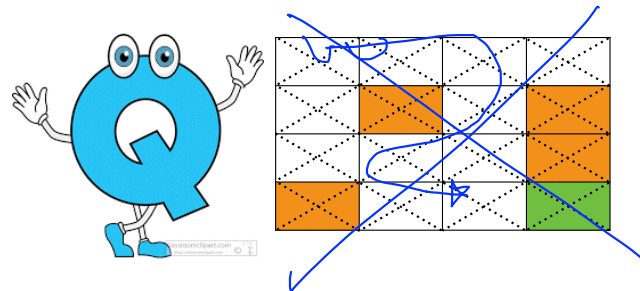
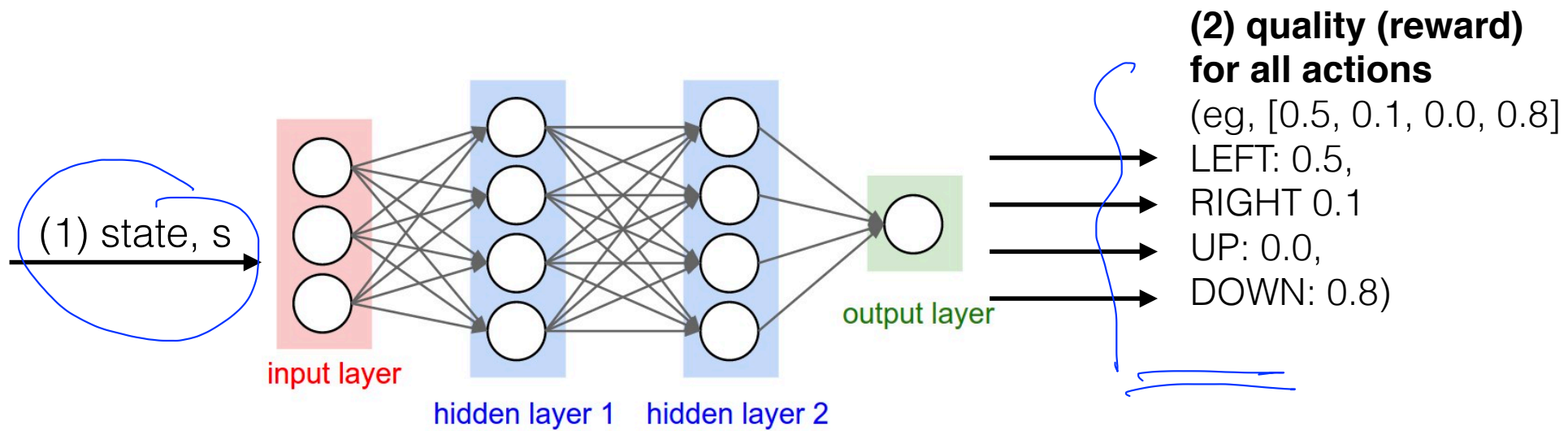


Lecture 7: DQN

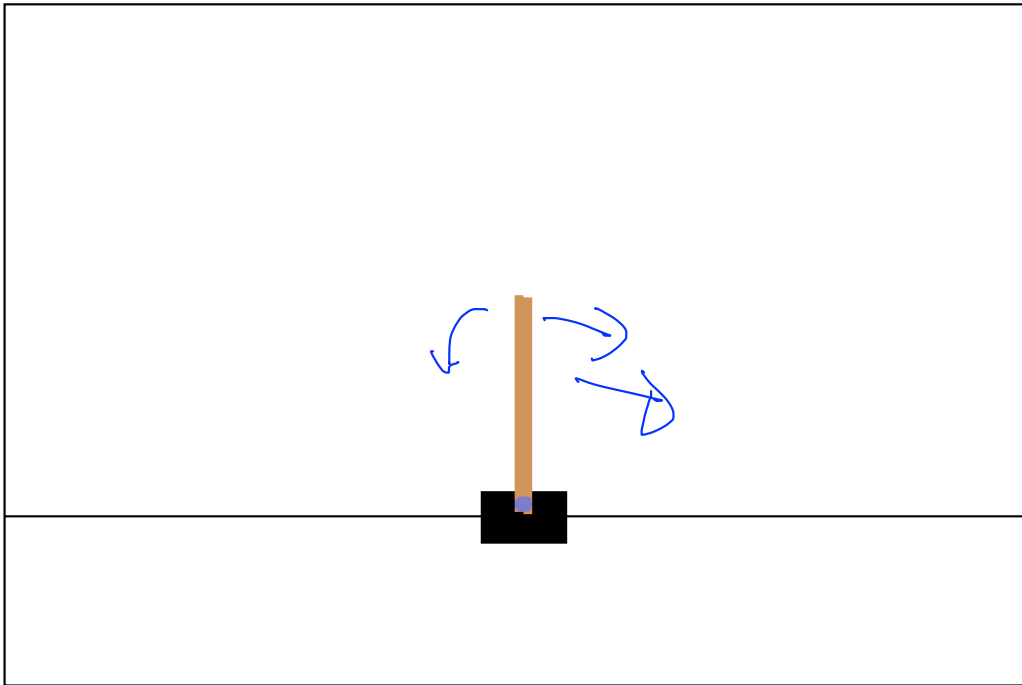
Reinforcement Learning with TensorFlow & OpenAI Gym
Sung Kim <hunkim+ml@gmail.com>



Q-function Approximation: Q-Nets



Q-Nets are unstable



```
Episode: 1988  steps: 14
Episode: 1989  steps: 25
Episode: 1990  steps: 15
Episode: 1991  steps: 23
Episode: 1992  steps: 19
Episode: 1993  steps: 17
Episode: 1994  steps: 46
Episode: 1995  steps: 20
Episode: 1996  steps: 17
Episode: 1997  steps: 15
Episode: 1998  steps: 33
Episode: 1999  steps: 22
2017-02-08 16:59:31.216 Pyth
Total score: 15.0
```

Convergence

\hat{Q} denote learner's current approximation to Q .

$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \theta))]^2$$

- ▶ Converges to Q^* using table lookup representation
- ▶ But **diverges** using neural networks due to:
 - ↳ Correlations between samples
 - ↳ Non-stationary targets

Reinforcement + Neural Net



There are some research papers on the topic:

- [Efficient Reinforcement Learning Through Evolving Neural Network Topologies \(2002\)](#)
- [Reinforcement Learning Using Neural Networks, with Applications to Motor Control](#)
- [Reinforcement Learning Neural Network To The Problem Of Autonomous Mobile Robot Obstacle Avoidance](#)

And some code:

- [Code examples](#) for neural network reinforcement learning.

Those are just some of the top google search results on the topic. The first couple of papers look like they're pretty good, although I haven't read them personally. I think you'll find even more information on neural networks with reinforcement learning if you do a quick search on Google Scholar.

- ▶ But **diverges** using neural networks due to:
 - ▶ Correlations between samples
 - ▶ Non-stationary targets

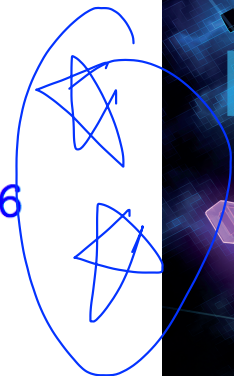
<http://stackoverflow.com/questions/10722064/training-a-neural-network-with-reinforcement-learning>

DQN paper

www.nature.com/articles/nature14236

DQN source code:

sites.google.com/a/deepmind.com/dqn/



Tutorial: Deep Reinforcement Learning, David Silver, Google DeepMind

Two big issues

- ▶ But **diverges** using neural networks due to:
 - ▶ Correlations between samples
 - ▶ Non-stationary targets

I. Correlations between samples

Algorithm 1 Deep Q-learning

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

end for

end for

I. Correlations between samples

Algorithm 1 Deep Q-learning

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

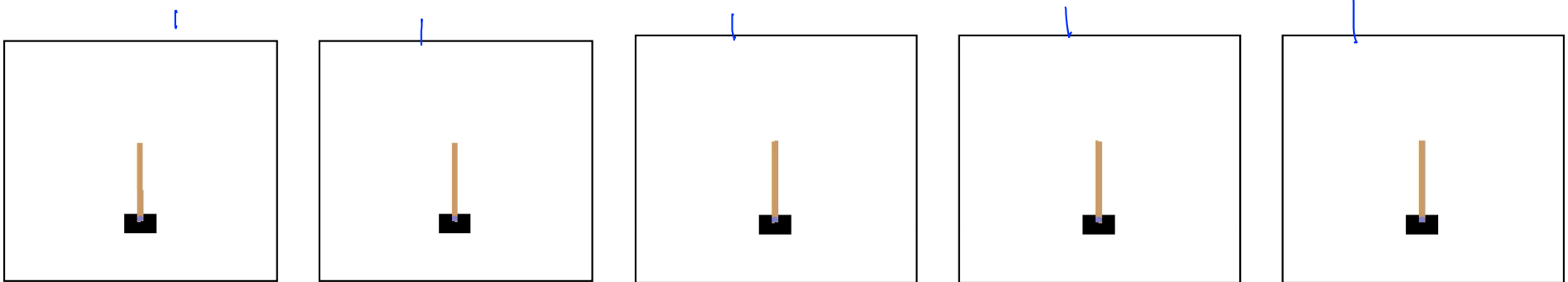


Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

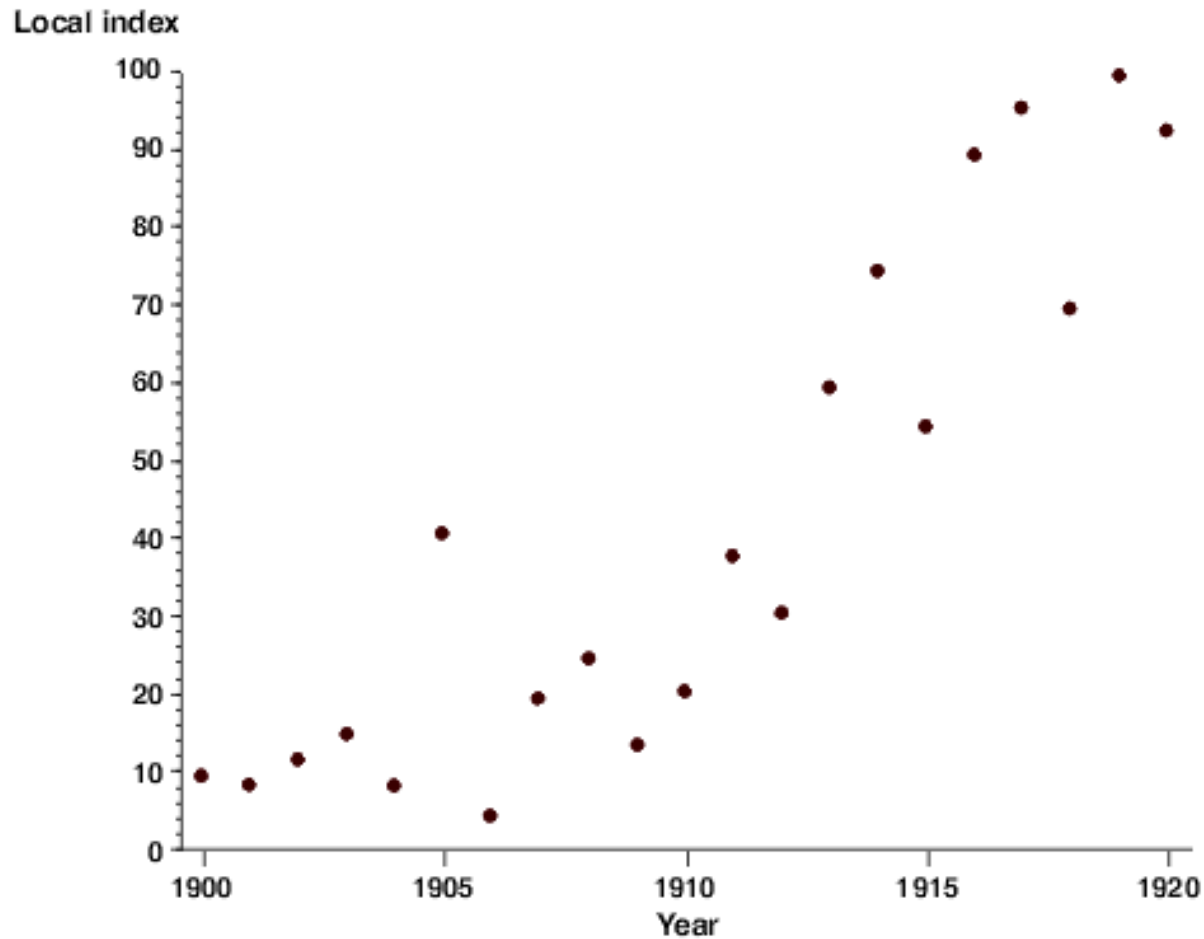
end for

end for



Playing Atari with Deep Reinforcement Learning - University of Toronto by V Mnih et al.

I. Correlations between samples



Prerequisite: <http://hunkim.github.io/ml/> or <https://www.inflern.com/course/기본적인-머신러닝-딥러닝-강좌/>

모두를 위한 딥러닝 강좌 시즌 1

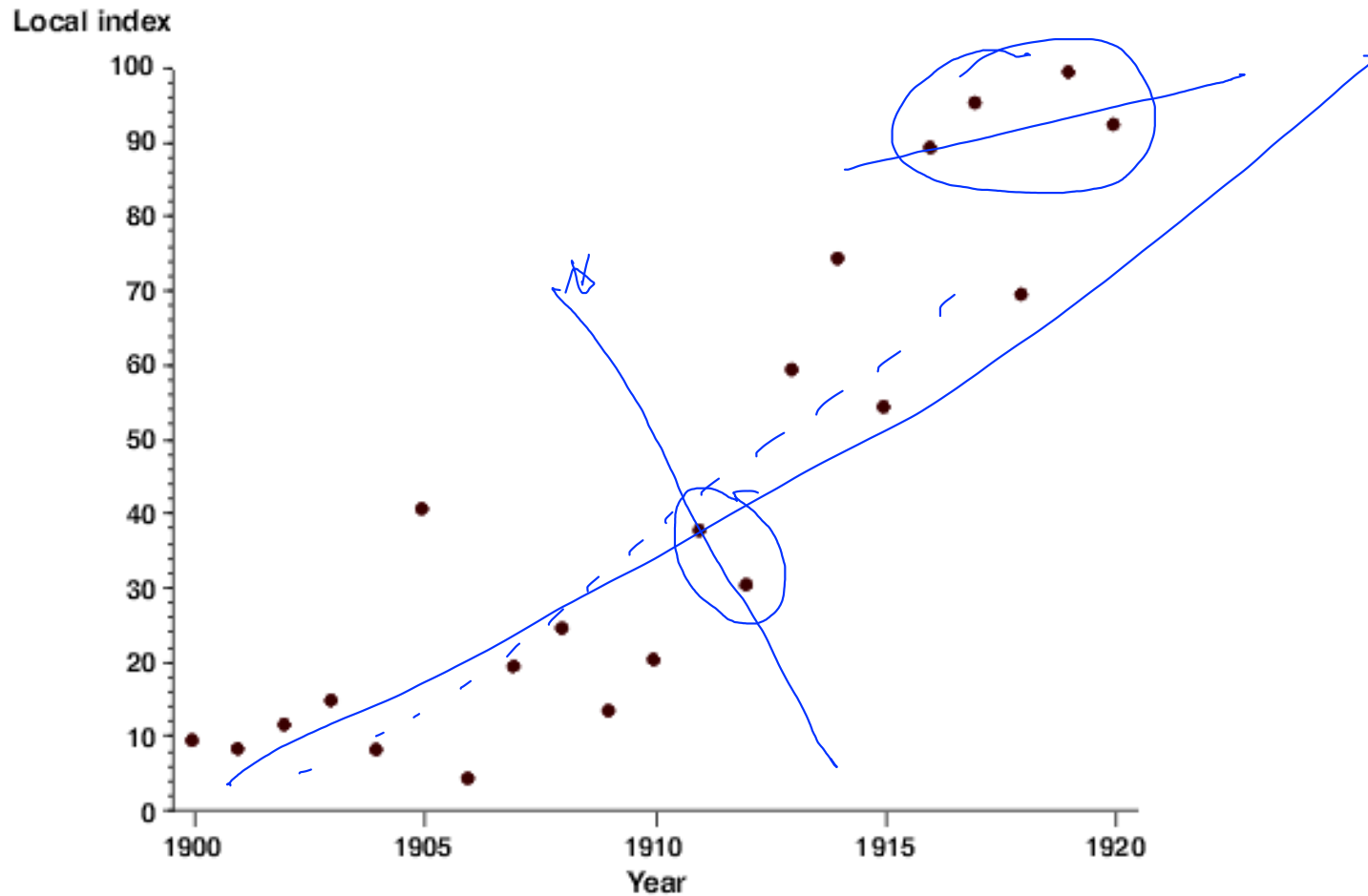
Sung Kim • 39 videos • 159,959 views • Updated today

Add a description

Play all Share Playlist settings Add videos

1		Lec 00 - Machine/Deep learning 수업의 개요와 일정 by Sung Kim	10:05
2		ML lec 01 - 기본적인 Machine Learnng의 용어와 개념 설명 by Sung Kim	14:27
3		ML lab 01 - TensorFlow의 설치 및 기본적인 operations by Sung Kim	10:48
4		ML lec 02 - Linear Regression의 Hypothesis 와 cost 설명 by Sung Kim	13:30
5		ML lab 02 - Tensorflow로 간단한 linear regression을 구현 by Sung Kim	10:00
6		ML lec 03 - Linear Regression의 cost 최소화 알고리즘의 원리 설명 by Sung Kim	16:12

I. Correlations between samples



2. Non-stationary targets

$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \theta))]^2$$



Q pred

\approx

\hat{Q}

target

2. Non-stationary targets

$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \theta))]^2$$

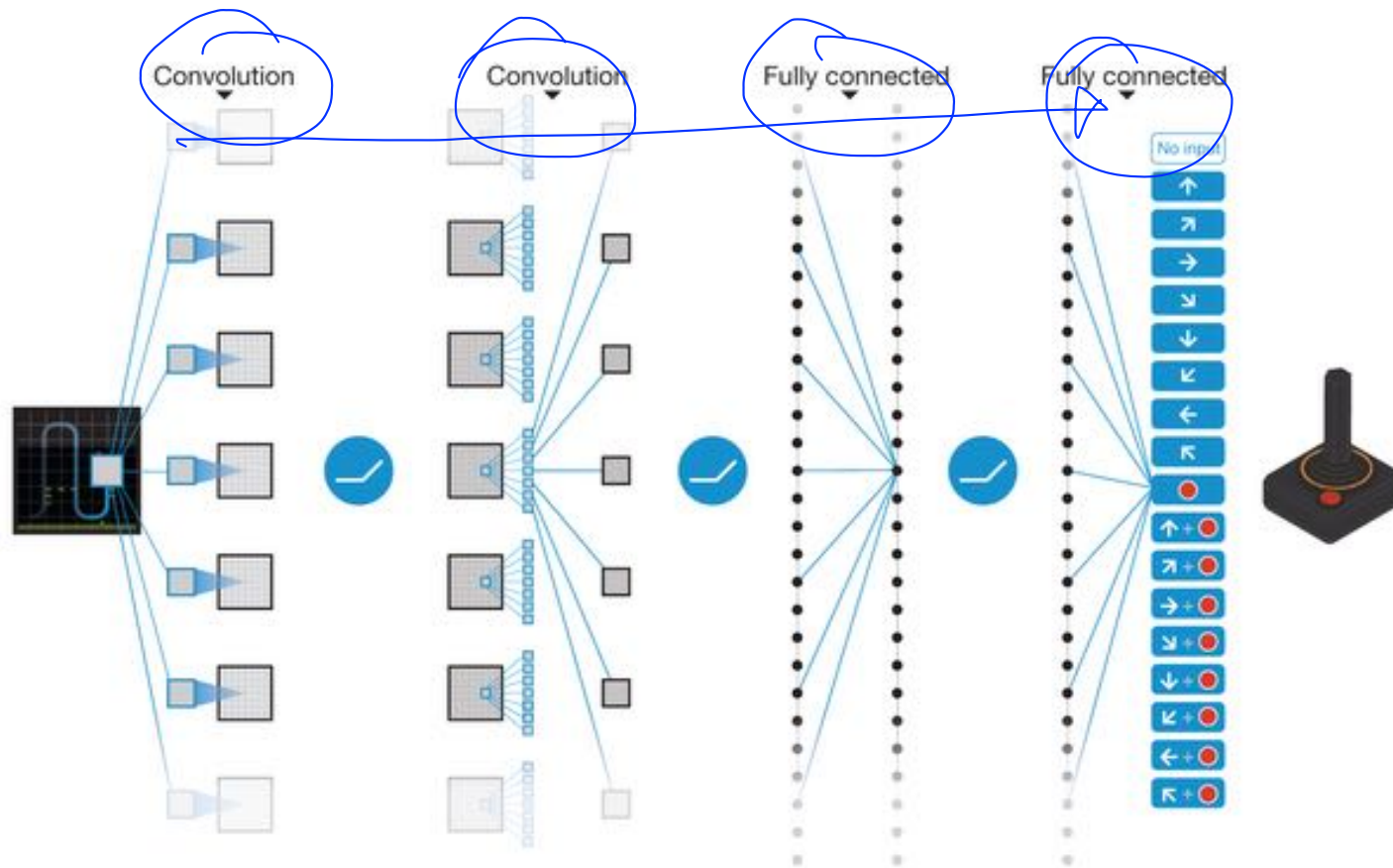
$\hat{Y} = \hat{Q}(s_t, a_t | \theta)$ $\hat{Y} = r_t + \gamma \max_{a'} \hat{Q}_{\theta}(s_{t+1}, a' | \theta)$

Handwritten annotations:
- Above the first \hat{Q} : pred
- Above the second \hat{Q} : target
- Below the first \hat{Q} : W
- Below the second \hat{Q} : W

DQN's three solutions

1. Go deep
2. Capture and replay
 - Correlations between samples
3. Separate networks: create a target network
 - Non-stationary targets

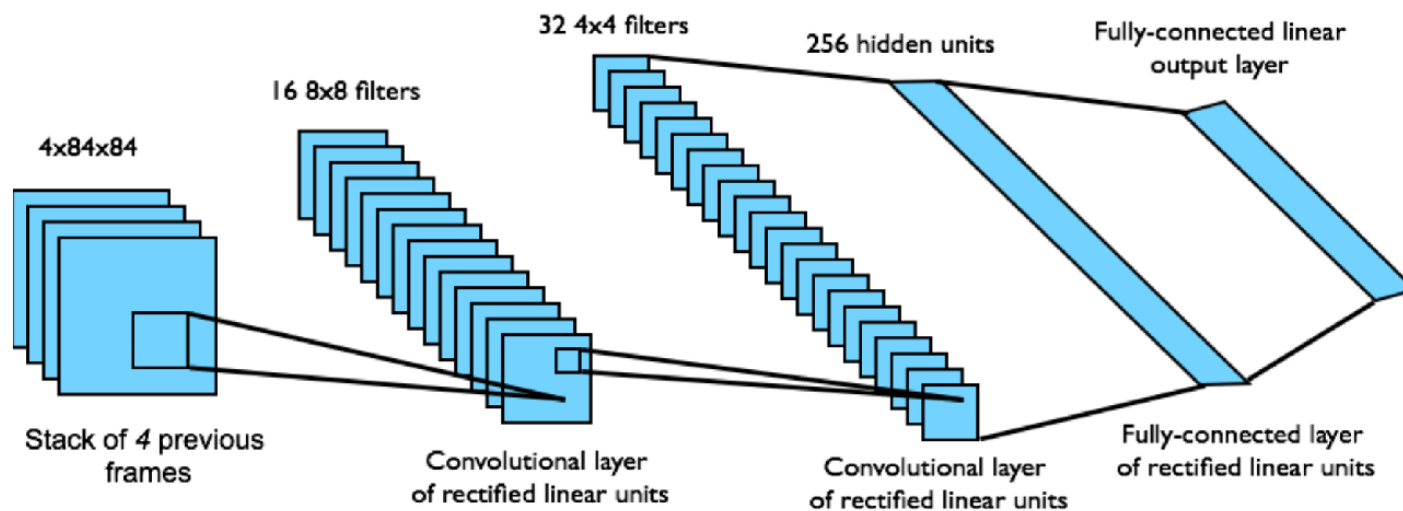
Solution 1: go deep



Human-level control through deep reinforcement learning, Nature
<http://www.nature.com/nature/journal/v518/n7540/full/nature14236.html>

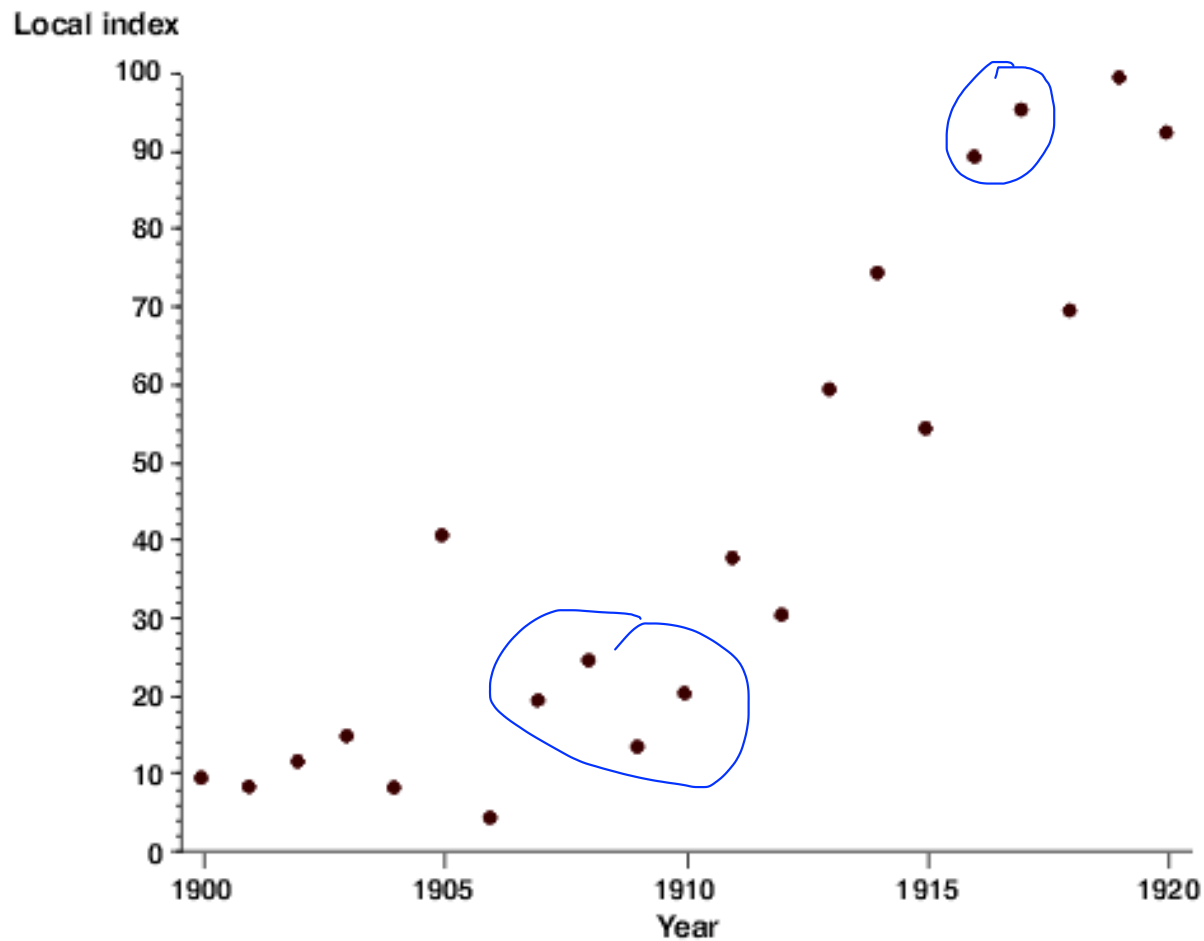
Solution I: go deep

- ▶ 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



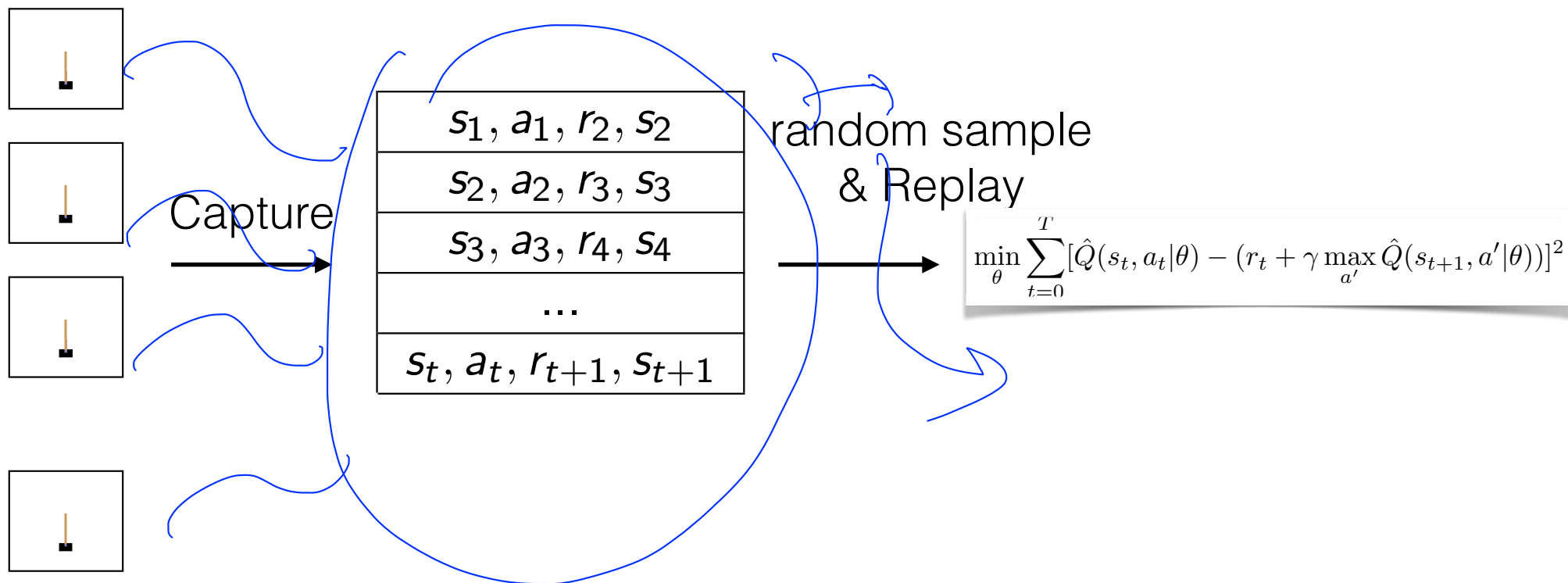
Network architecture and hyperparameters fixed across all games

Problem 2: correlations between samples



2019

Solution 2: experience replay



Solution 2: experience replay

Algorithm 1 Deep Q-learning with Experience Replay

Initialize action-value function Q with random weights

for episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

for $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in \mathcal{D}

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from \mathcal{D}

 Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

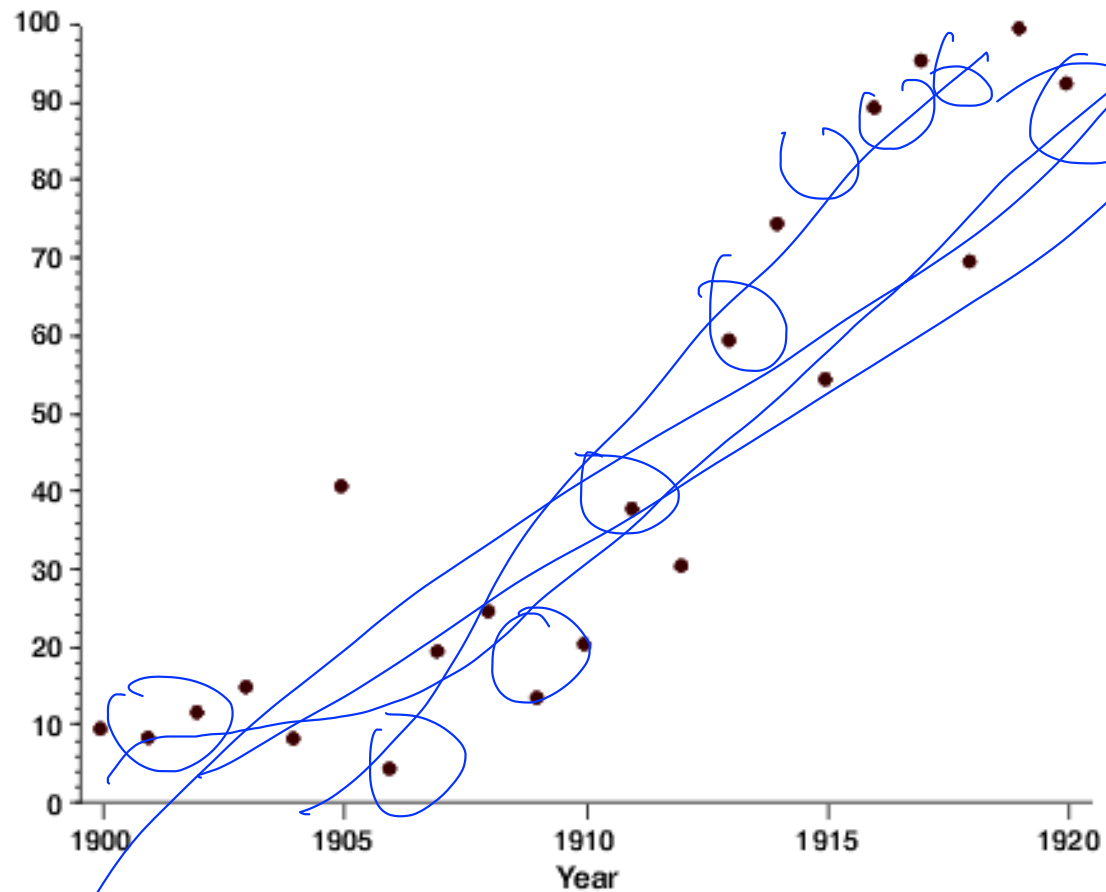
end for

end for

Problem 2: correlations between samples

s_1, a_1, r_2, s_2
s_2, a_2, r_3, s_3
s_3, a_3, r_4, s_4
...
$s_t, a_t, r_{t+1}, s_{t+1}$

Local index



Problem 3: non-stationary targets

$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \theta))]^2$$

$$\hat{Y} = \hat{Q}(s_t, a_t | \theta)$$

Handwritten annotations: A blue arrow points from the \hat{Q} term in the equation above to the \hat{Y} term. Another blue arrow points from the θ parameter in the equation above to the θ parameter in the equation below.

$$Y = r_t + \gamma \max_{a'} \hat{Q}_{\theta}(s_{t+1}, a' | \theta)$$

Handwritten annotation: The word "long" is written in blue above the $\max_{a'}$ term.

Solution 3: separate target network

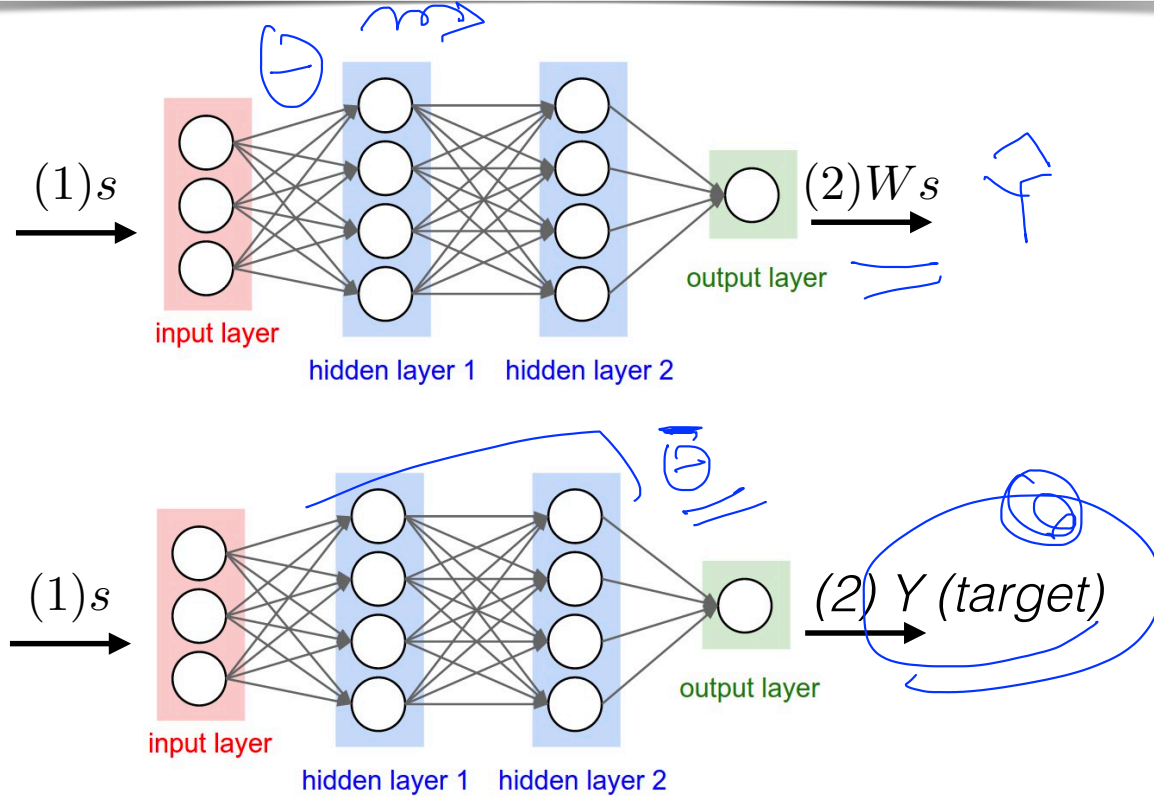
$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \bar{\theta}))]^2$$

~~$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \theta))]^2$$~~

Human-level control through deep reinforcement learning, Nature
<http://www.nature.com/nature/journal/v518/n7540/full/nature14236.html>

Solution 3: separate target network

$$\min_{\theta} \sum_{t=0}^T [\hat{Q}(s_t, a_t | \theta) - (r_t + \gamma \max_{a'} \hat{Q}(s_{t+1}, a' | \bar{\theta}))]^2$$



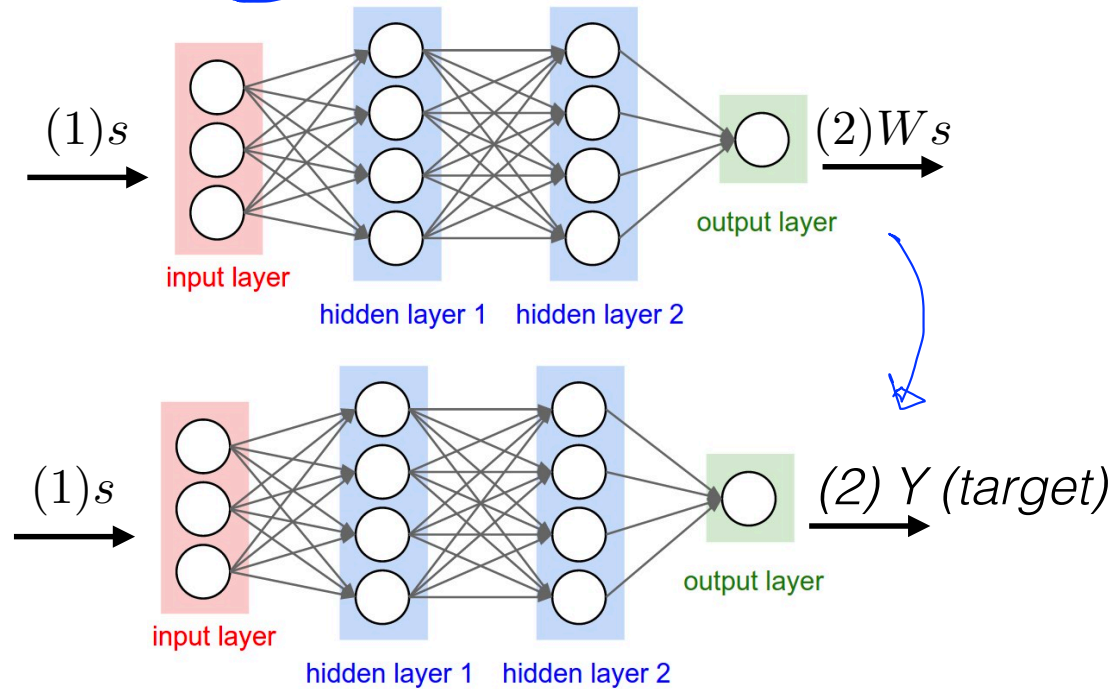
Solution 3: copy network

Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from D

Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

Perform a gradient descent step on $(y_j - Q(\phi_j, a_j, \theta))^2$ with respect to the network parameters θ

Every C steps reset $\hat{Q} = Q$



Understanding Nature Paper (2015)

Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ

Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$

For episode = 1, M do

Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

For $t = 1, T$ do

With probability ϵ select a random action a_t

otherwise select $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

Execute action a_t in emulator and observe reward r_t and image x_{t+1}

Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D

Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from D

Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a; \theta^-) & \text{otherwise} \end{cases}$

Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters θ

Every C steps reset $\hat{Q} = Q$

End For

End For

Next
Lab: DQN

