

Lecture 6: Q-Network

Reinforcement Learning with TensorFlow&OpenAI Gym

Sung Kim <hunkim+ml@gmail.com>

Q-Table (16x4)

(1) state, s →
(2) action, a →



(3) quality (reward)
for the given action
(eg, LEFT: 0.5, RIGHT 0.1
UP: 0.0, DOWN: 0.8)
→

$Q(s, a)$

Q-learning Test

created by Jae Hyun Lee(jaehyunlee25@gmail.com)

STATUS

EPISODE: 0 | MOVE: 0 times | TIME ELAPSED: 0.0 secs | Hole: 0 times | GOAL: 0 times

OPTIONS

SIZE: 4 | VELOCITY: normal | REPEAT: 1000 * 1 (1000) times

Hole: reward -1 penalty go on | Gamma(γ): 0.9 | Explo it

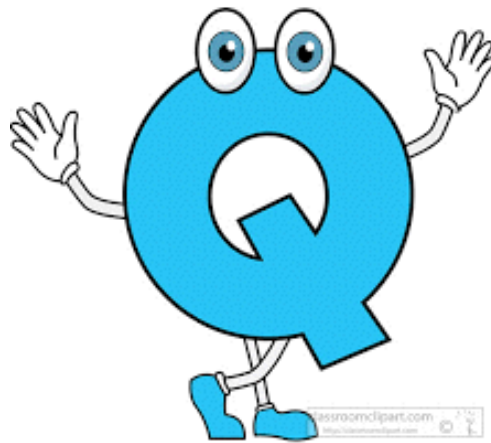
START

0 0 0 0 R=0	0 0 0 0 R=0	0 0 0 0 R=0	0 0 0 0 R=0
0 0 0 0 R=0	0 0 0 0 R=-1	0 0 0 0 R=0	0 0 0 0 R=-1
0 0 0 0 R=0	0 0 0 0 R=0	0 0 0 0 R=0	0 0 0 0 R=-1
0 0 0 0 R=-1	0 0 0 0 R=0	0 0 0 0 R=0	0 0 0 0 R=1

<http://computingkoreanlab.com/app/jAI/jQLearning/>

Q-Table (16x4)

(1) state, s
→
(2) action, a
→



(3) quality (reward)
for the given action
(eg, LEFT: 0.5, RIGHT 0.1
UP: 0.0, DOWN: 0.8)
→

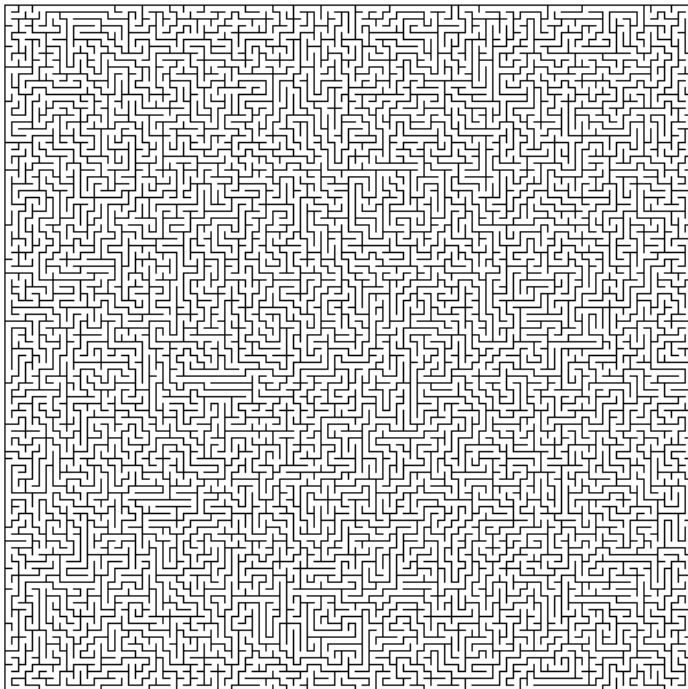
$Q(s, a)$



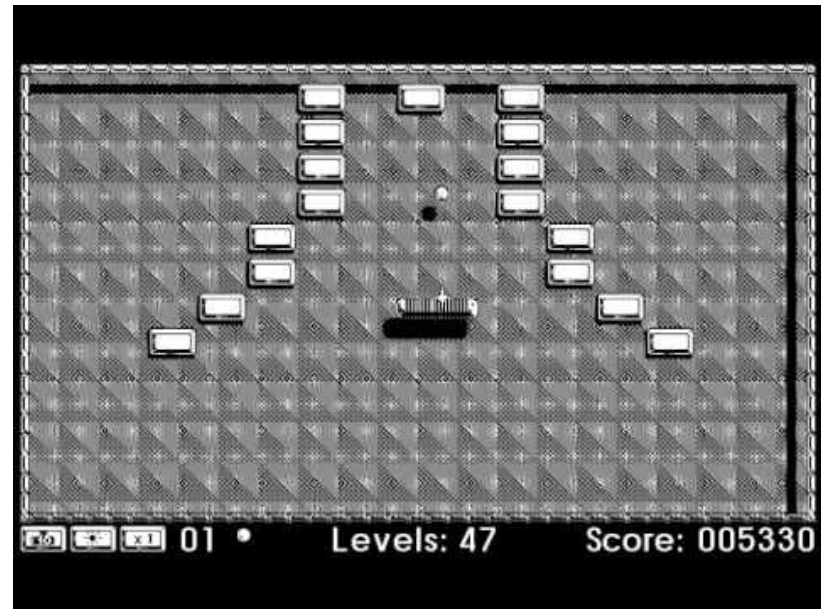
Q-Table (?)

160x160x4

2 80x80



100x100 maze



80x80 pixel + 2 color (black/white)

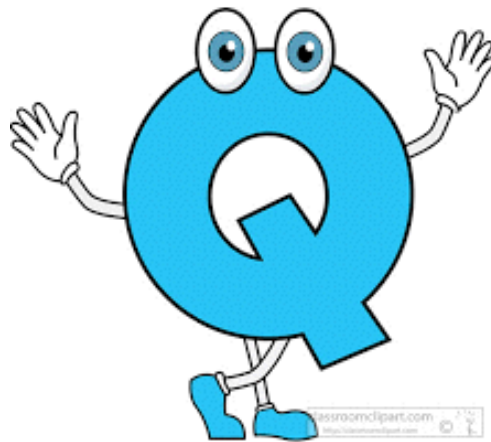


pow(2, 80*80)?

```
Sungs-MacBook-Pro:qlearning hunkim$ python
Python 2.7.10 (default, Jul 30 2016, 19:40:32)
[GCC 4.2.1 Compatible Apple LLVM 8.0.0 (clang-800.0.34)] on darwin
Type "help", "copyright", "credits" or "license" for more information.
>>> pow(2, 80*80)
390815922664323873317461428361483673112676810704634122725066747207685355684301383817204958869117433070781824
345011844830281272995124732157651316662436942651956617552114606076508167997675804172067930054414257889999968
221814559073171121585237586191625968426486695394453385378020623210968986095655234800673206169578993893064621
394479007844522654750932530260932906944511708573961116842051110072780702996021975530468334392852283568123658
054474934577299768778927200578050584569281027959847481492880157513920578725804851336909371096131251468207589
952539391748619172204401099253655425443334675779740145331744115766832347511737830030067538031141178372089229
186020884093542742187687849517214364478837908382646195523281445267002410663402978244431485772594698434066258
955213768118705885537084067810415837310291314293532224816692313560648857707617678657107208778727572115167144
969844554782437658117922884282324307657985008269944091087969017270165493533858541923739452891021946640039871
139014694517482748438203362034911702773985087028549861930279639065901103098389957112947553195190069994199572
871577997339201518554694093470424633134190567138126513775817867770518591303000260334804048037845284233493883
534489638498647260587064326521948900609819006765410832117094339235540061356424944564563951032375578180928992
076476796598306704845779425650491079898232492863540154374240499247068529525671271000571306646256947045781135
744092881405282687148040508264376853441383821755805287956740468541933707091945416554391392625479753506217540
329184028243756578451089052788428260969378835272884507542718466673627090471859996737723112934060270454109059
566034439294502155993552543831988709035361713548867020994349284913996584689674031362649588710526967617555709
716501891685314826079439170843819922088878128902968582950531577390213885399071760414288571960169709472040532
8129745119056694810317474513400330333517233611193133379954007384384395034058799871873253376L
>>>
```


Q-Table (? x ?)

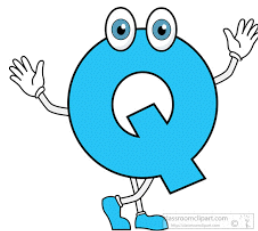
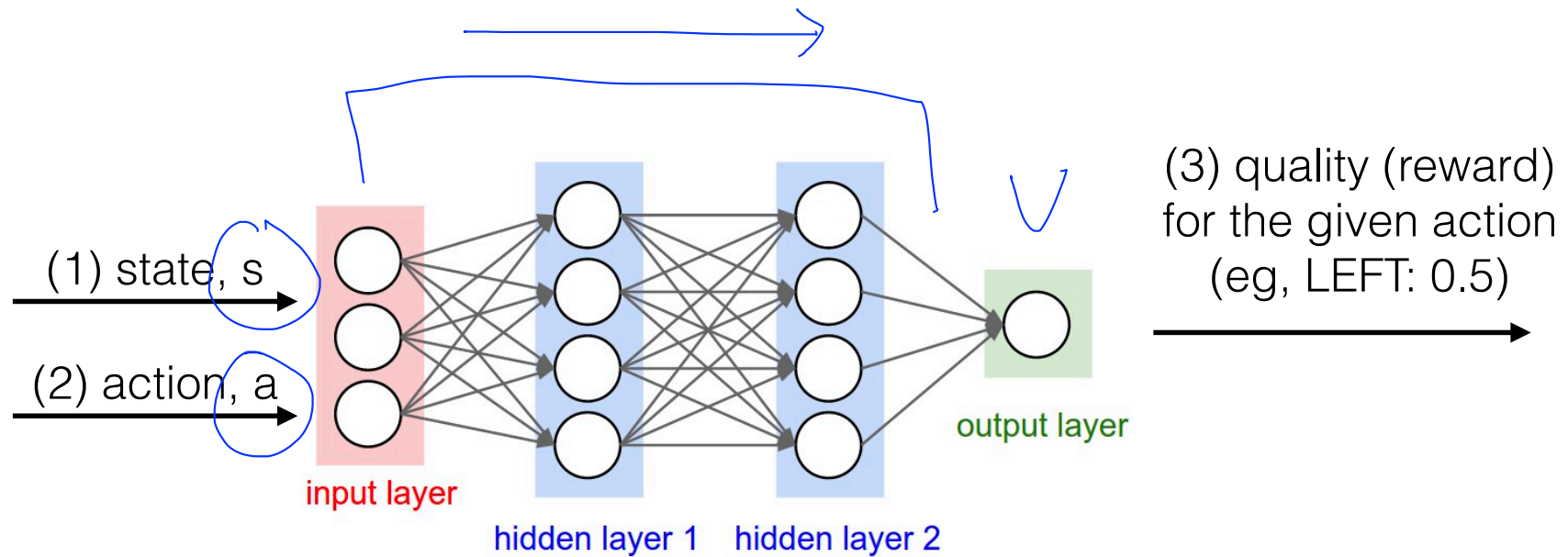
(1) state, s
→
(2) action, a
→



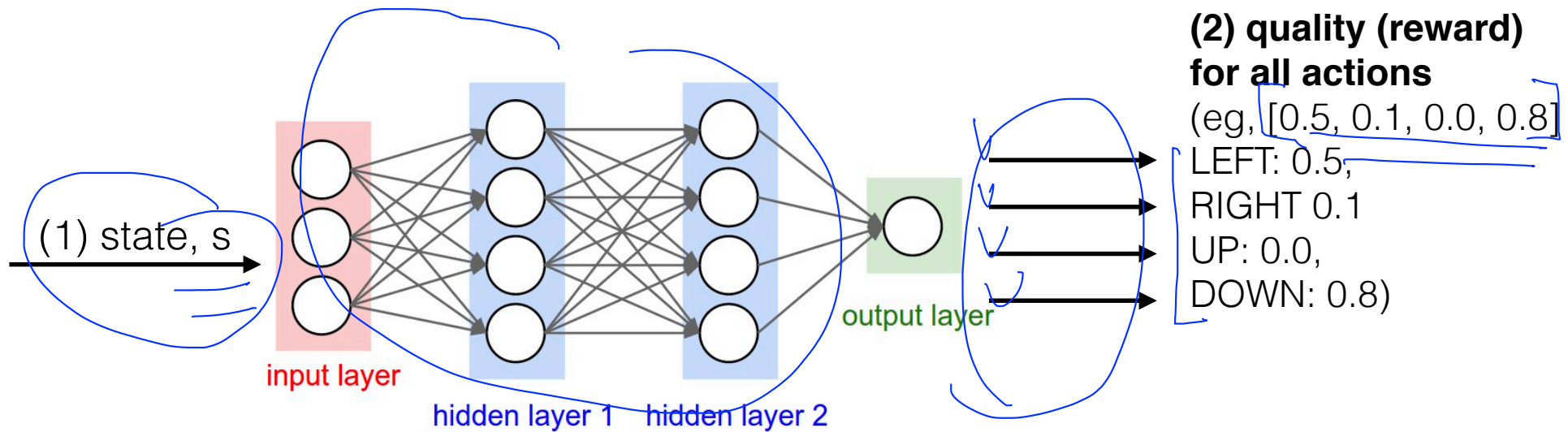
$Q(s, a)$

(3) quality (reward)
for the given action
(eg, LEFT: 0.5)
→

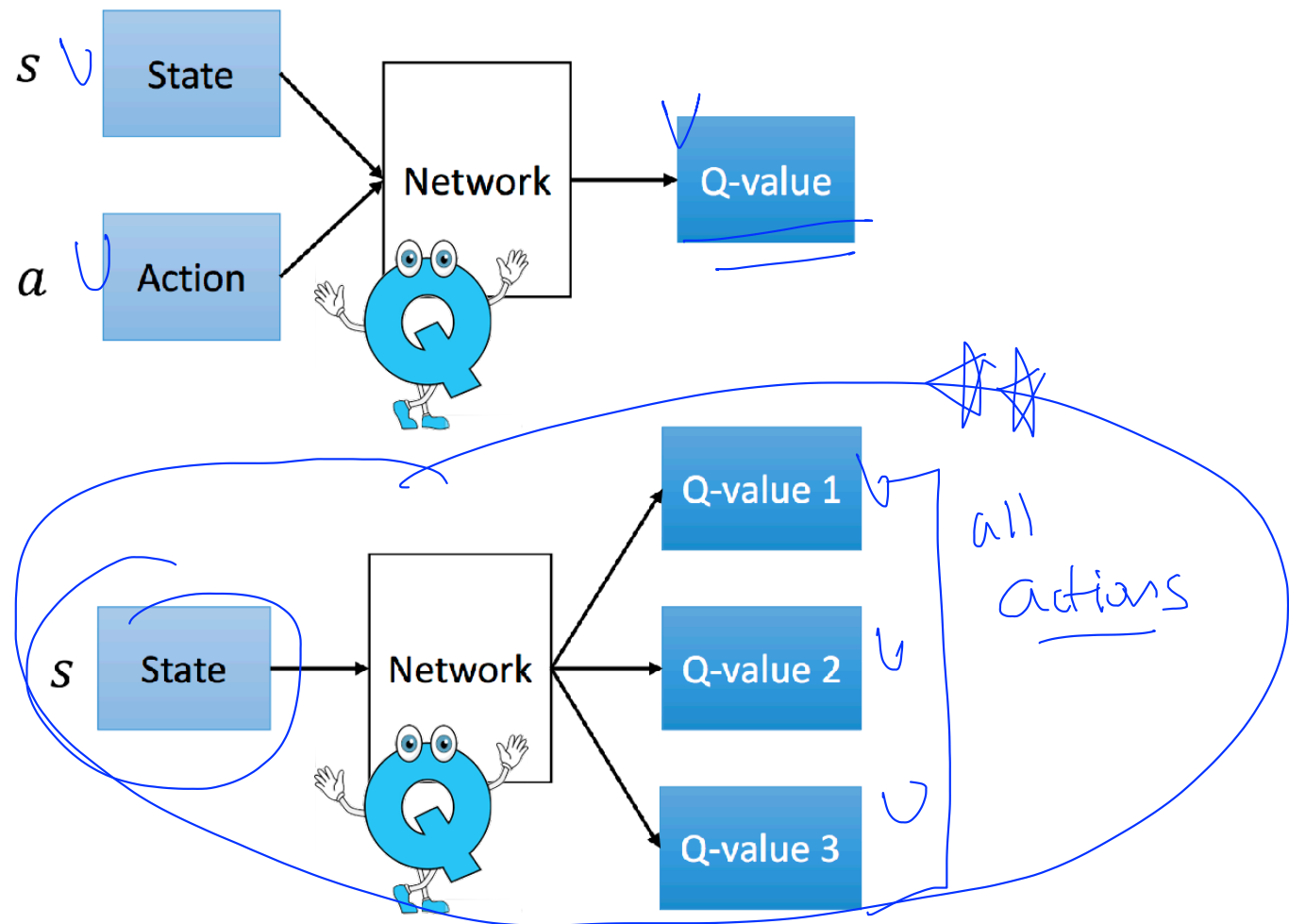
Q-function Approximation



Q-function Approximation



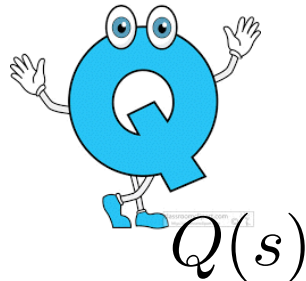
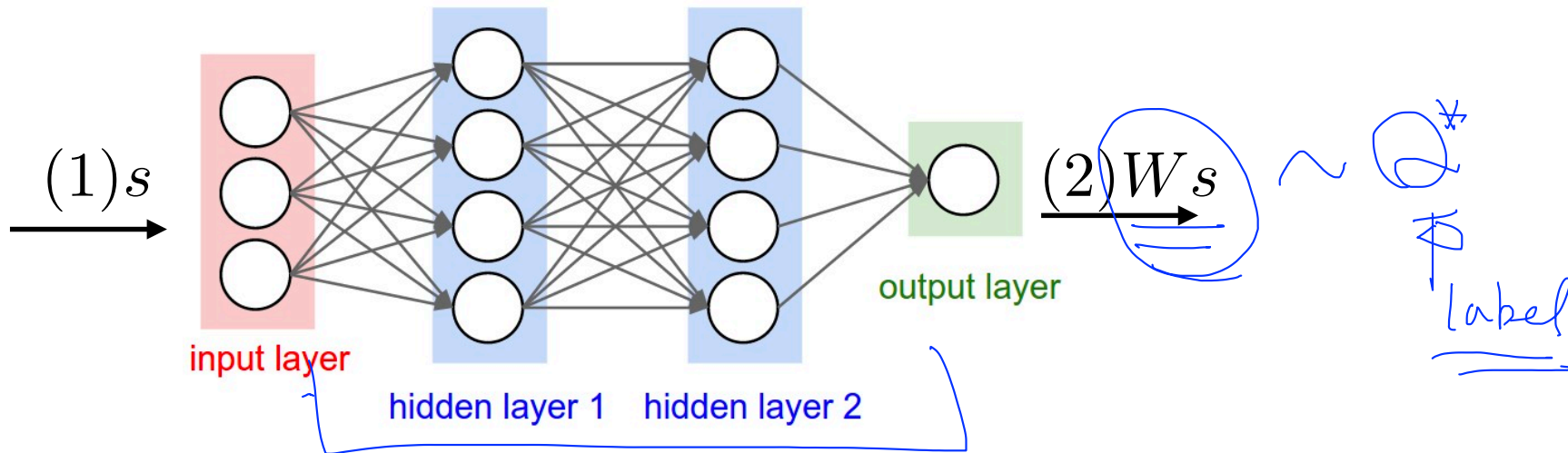
Q-function Approximation




Q-Network training (linear regression)


$$H(\underline{x}) = \underline{W}x$$


$$\underline{\text{cost}}(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2$$





Prerequisite: <http://hunkim.github.io/ml/> or <https://www.infllearn.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 Learning의 용어와 개념 설명
by Sung Kim More X

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

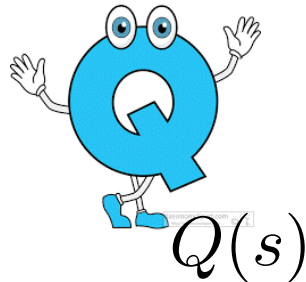
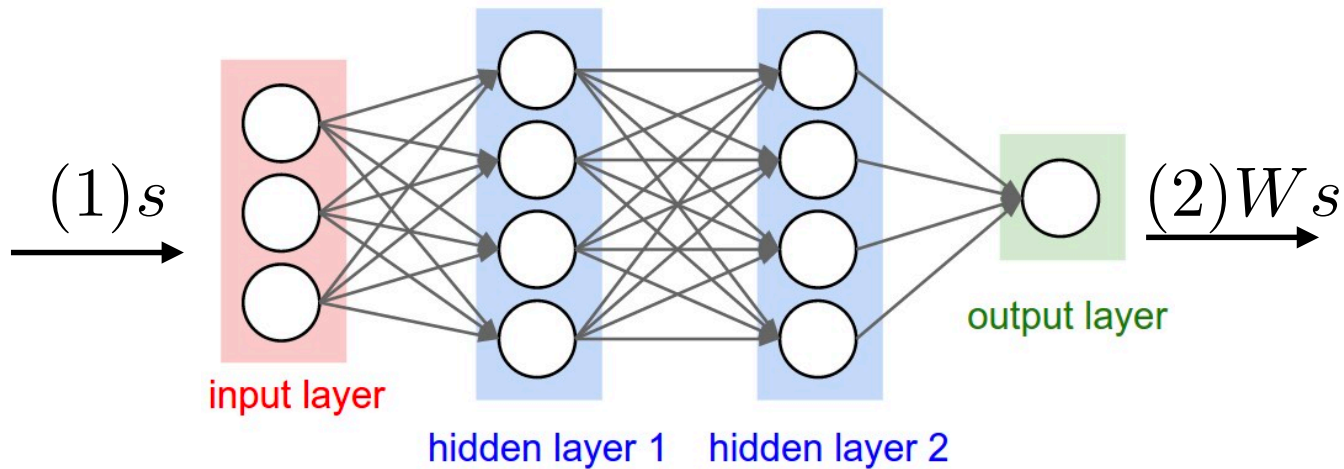
5  ML lab 02 - Tensorflow로 간단한 linear regression을 구현
by Sung Kim 10:00

6  ML lec 03 - Linear Regression의 cost 최소화 알고리즘의 원리 설명
by Sung Kim 16:12

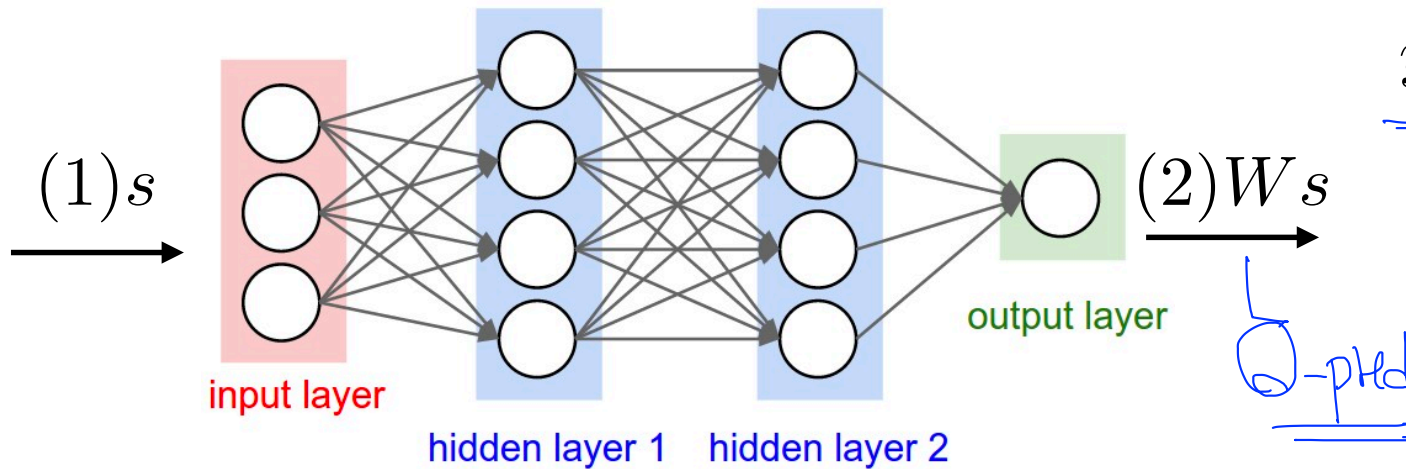
Q-Network training (linear regression)

$$\underline{H(x) = Wx}$$

$$\underline{cost(W) = \frac{1}{m} \sum_{i=1}^m (Wx^{(i)} - y^{(i)})^2}$$



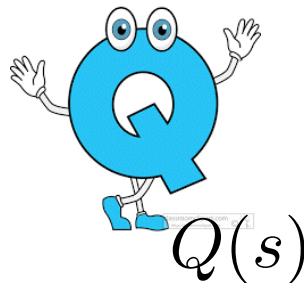
Q-Network training (linear regression)



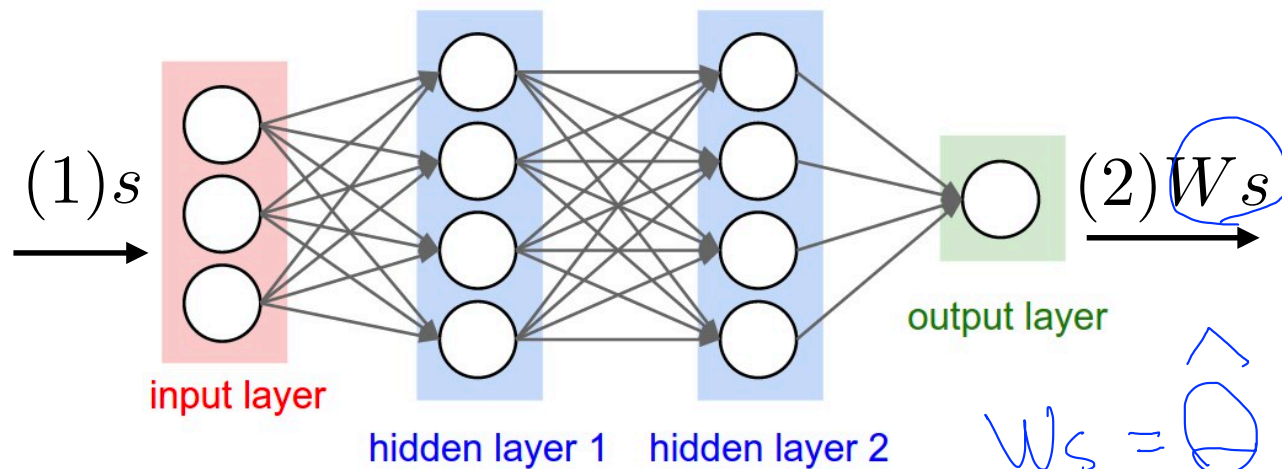
$$\text{cost}(W) = (\underline{W} s - \underline{y})^2$$

$$\underline{y} = r + \gamma \max_{Q^*} Q(s')$$

Q -pred



Q-Network training (math notations)



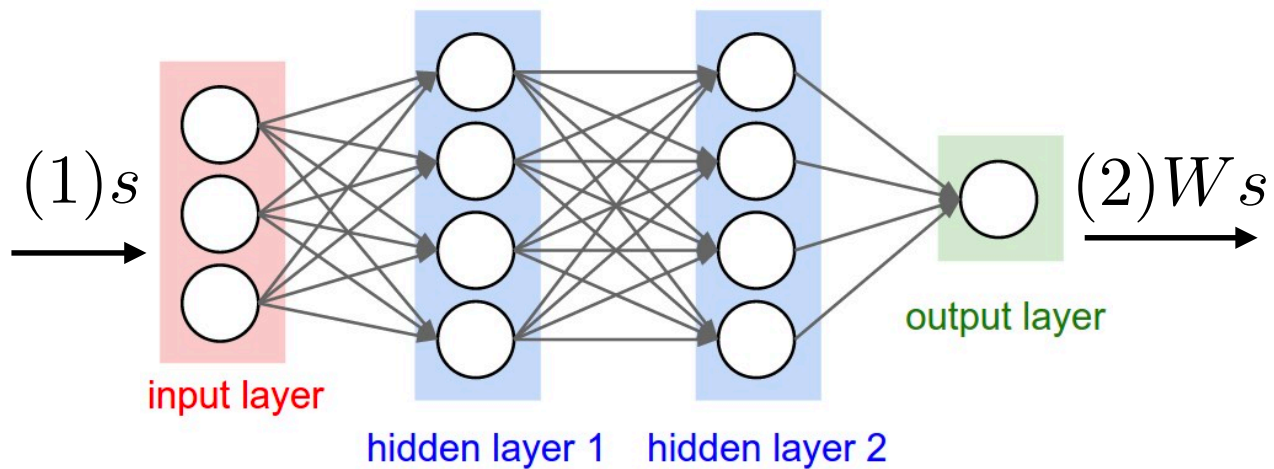
- Approximate Q^* function using θ
- $$\hat{Q}(s, a | \theta) \sim \underline{Q^*}(s, a)$$

$$W_s = \hat{Q}(s, a | \theta)$$

weight
(Let work)



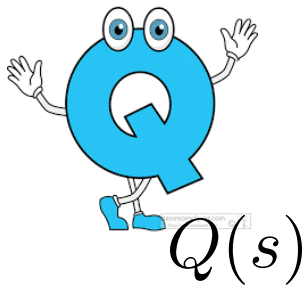
Q-Network training (math notations)



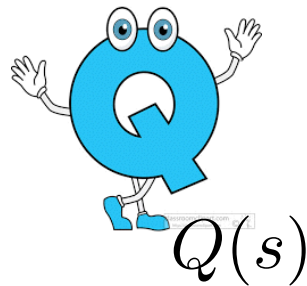
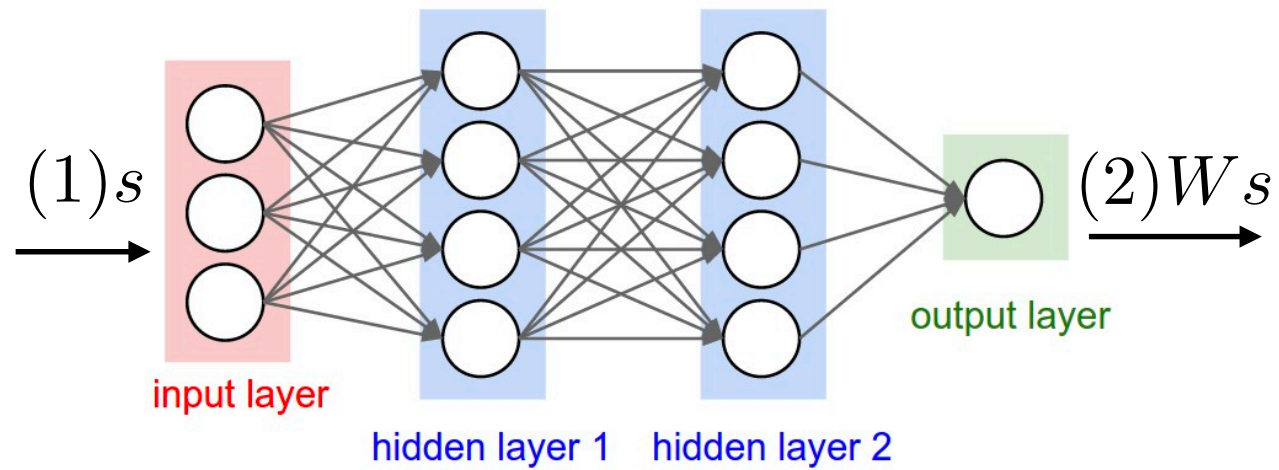
- Approximate Q^* function using θ
 $\hat{Q}(s, a|\theta) \sim Q^*(s, a)$

- Choose θ to minimize

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



Q-Network



Algorithm

$$\leftarrow \phi(s) \quad \phi \approx s$$

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 $\phi_{t+1} = \phi(s_{t+1})$

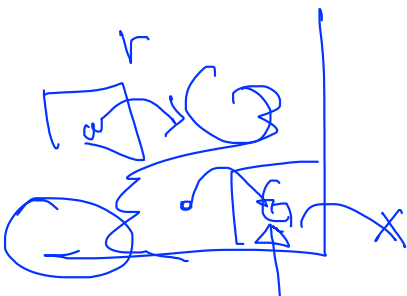
Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

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

end for

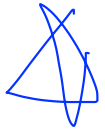
end for

Y label and loss function



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

Deterministic or Stochastic?



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

Handwritten notes: "target?" with an arrow pointing to the max term, and a box around the entire set definition.

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

Handwritten notes: "target" with an arrow pointing to the first Q term, and a large box around the equation.

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

Handwritten note: "Q" with an arrow pointing to the Q term in the max.

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:

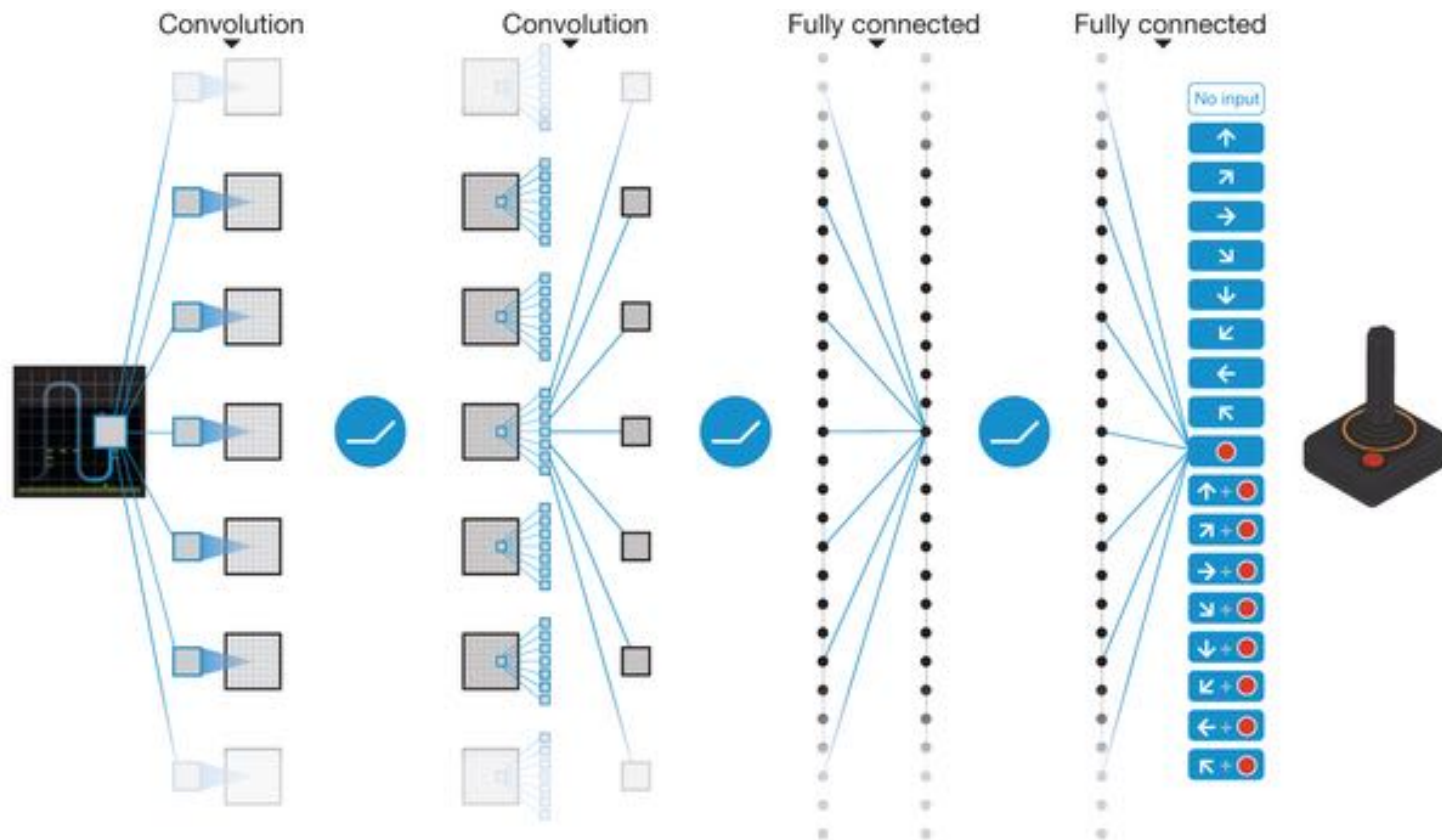
- [Code examples](#) for neural network reinforcement learning.

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: Deep, Replay, Separated networks



Next
Lab: Q-network

