



# Practical Machine Learning

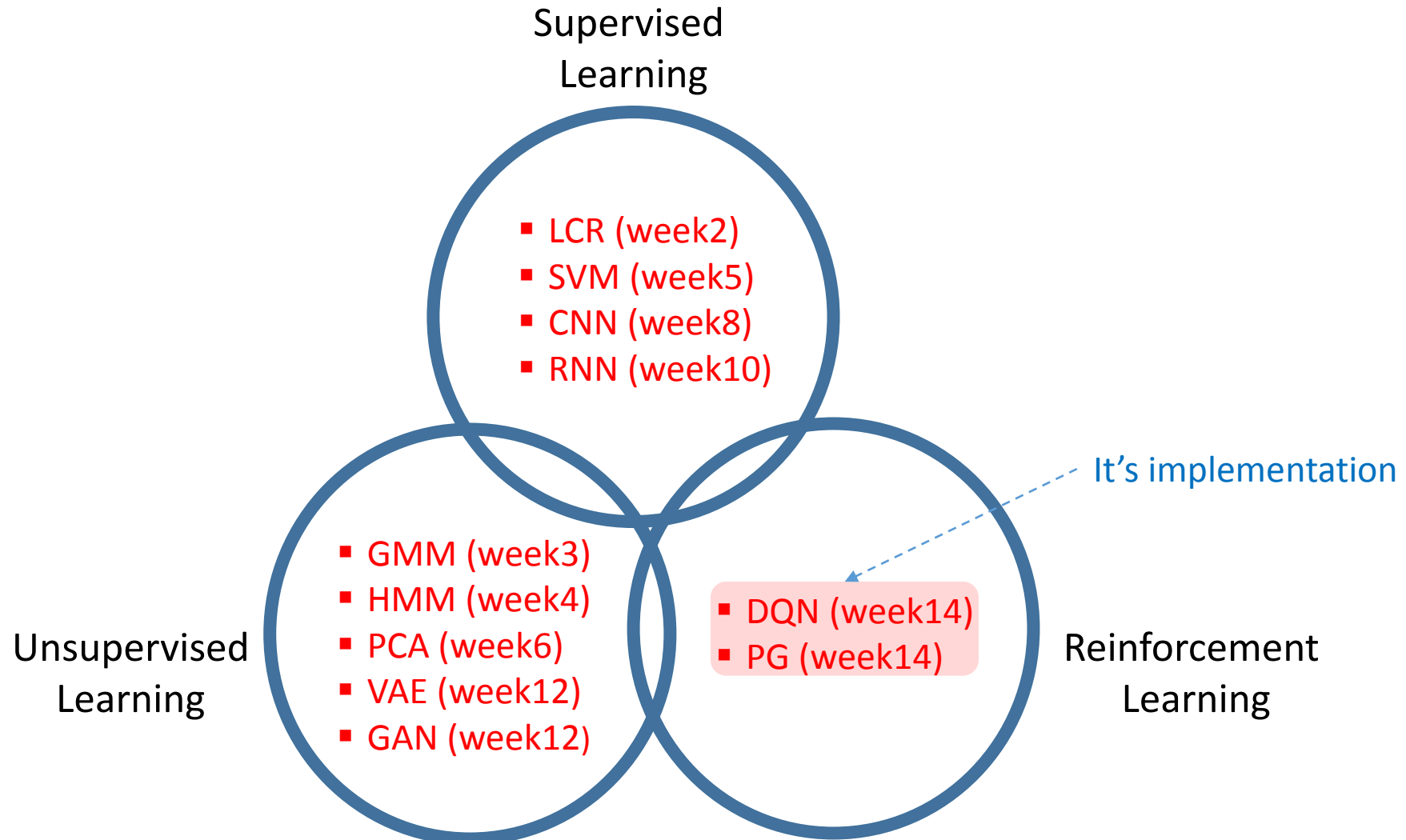
## Lecture 15

### Tensorflow – DQN/PG/AC implementation

Dr. Suyong Eum



# Where we are



# You are going to learn


- ❑ What OpenAI and Gym are,
- ❑ Implementations of
  - Deep Q-Network (DQN) – 2015
  - Policy Gradient (PG)
  - Actor Critic (AC)

## ❑ OpenAI :

- A non-profit artificial intelligence (AI) research company that aims to promote and develop friendly AI in such a way to benefit humanity as a whole.
- In October 2015, Elon Musk et al founded the organization.
- On April 2016, OpenAI released a public beta of “OpenAI Gym”, its platform for reinforcement learning research.

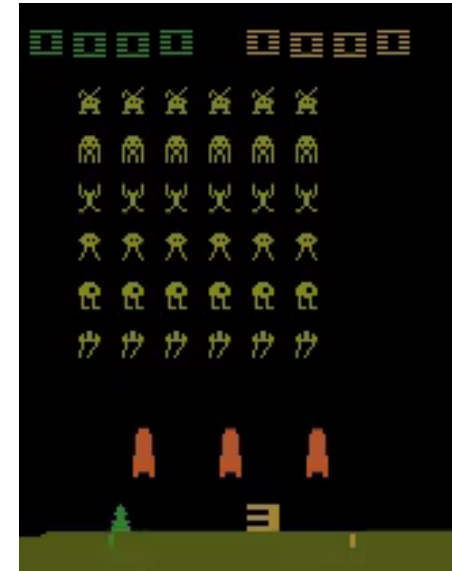
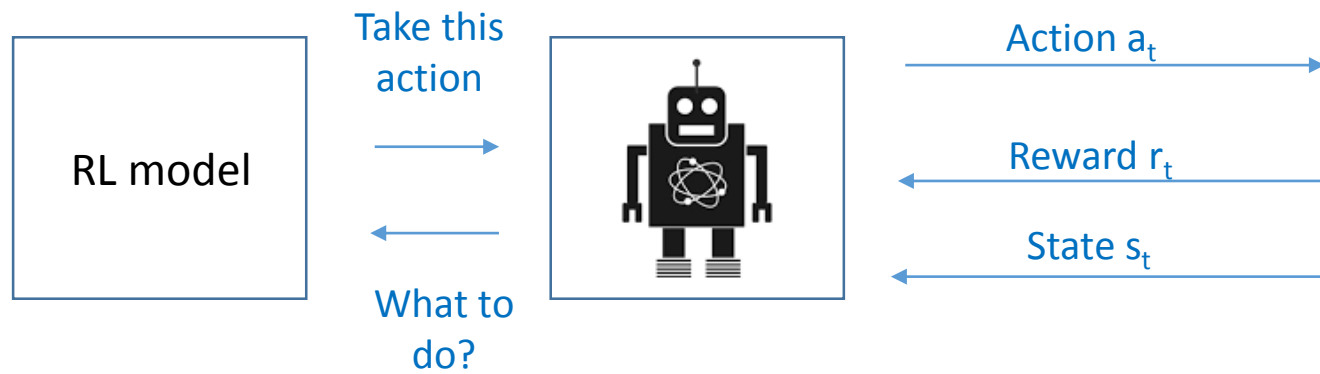
## ❑ OpenAI Gym

- A toolkit for developing and comparing reinforcement learning algorithms
- <https://github.com/openai/gym>
- <https://gym.openai.com/>

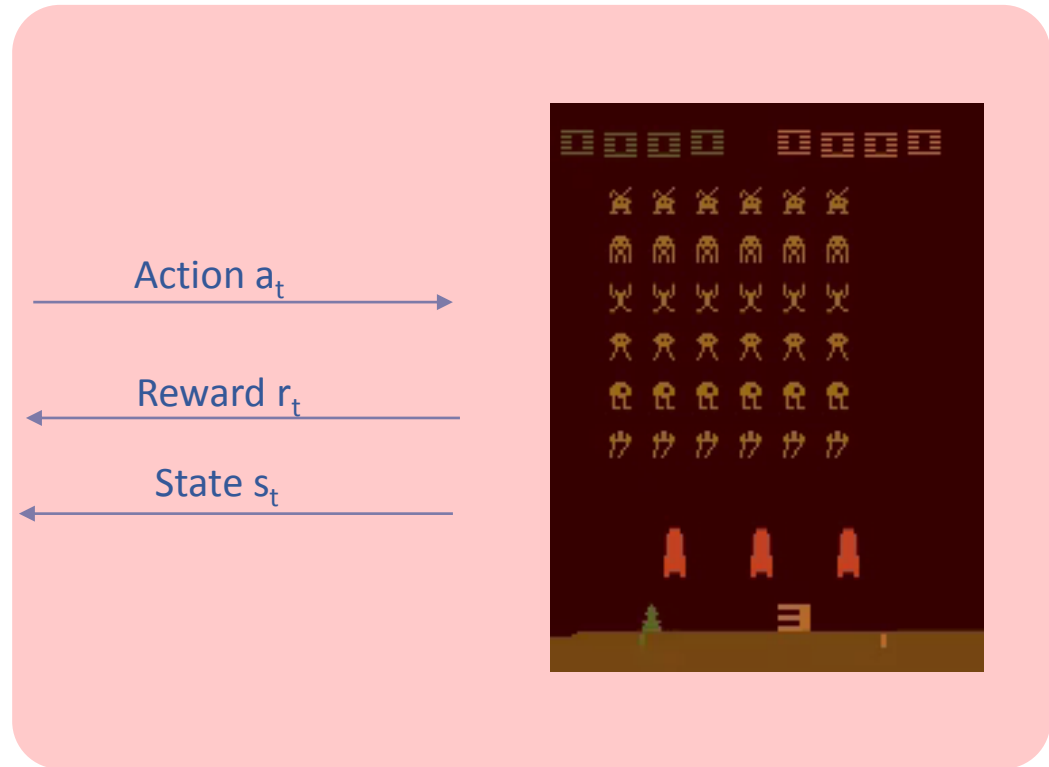
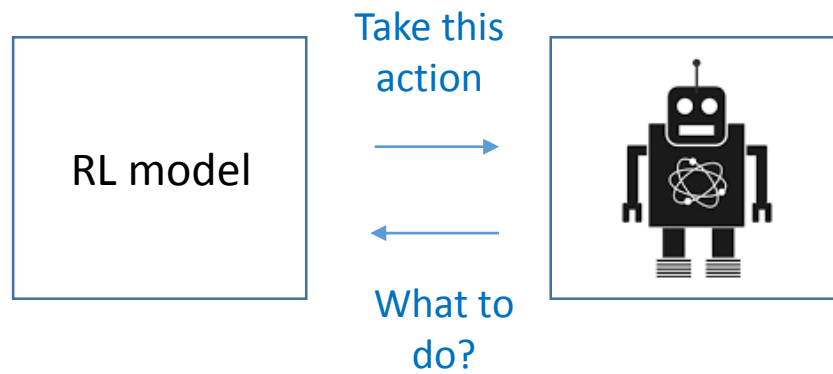


Gym

Gym is a toolkit for developing and comparing reinforcement learning algorithms. It supports teaching agents everything from walking to playing games like Pong or Pinball.

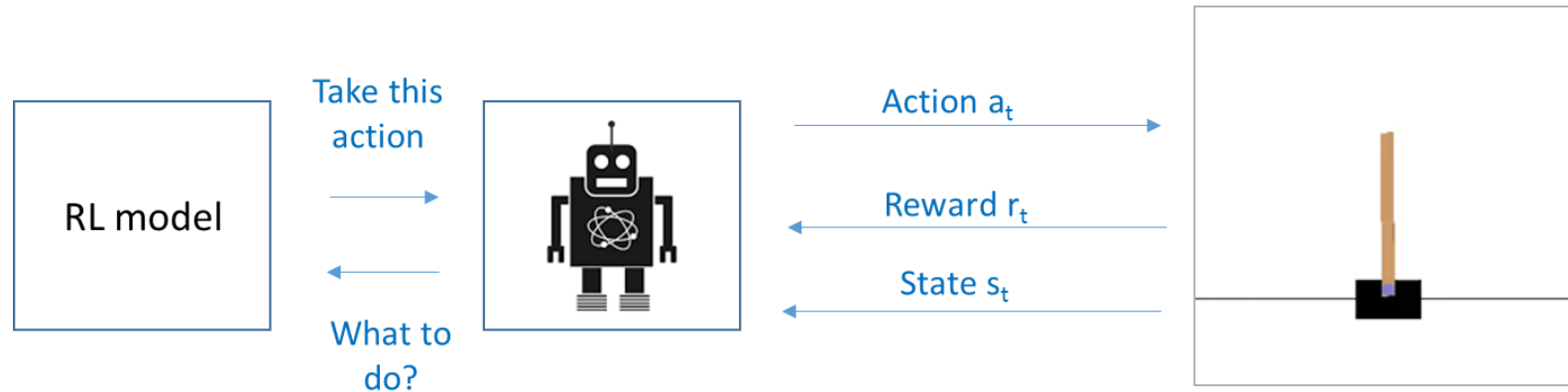


# OpenAI Gym framework



Gym framework

# OpenAI Gym framework: CartPole game



```
import gym
env=gym.make('CartPole-v0')  # Which game you want to play?

# Initialize the current state
current_state = env.reset()  # Initialization

# Action as an input to gym
action = 1 # right  # Two actions: right(1) or left (0)

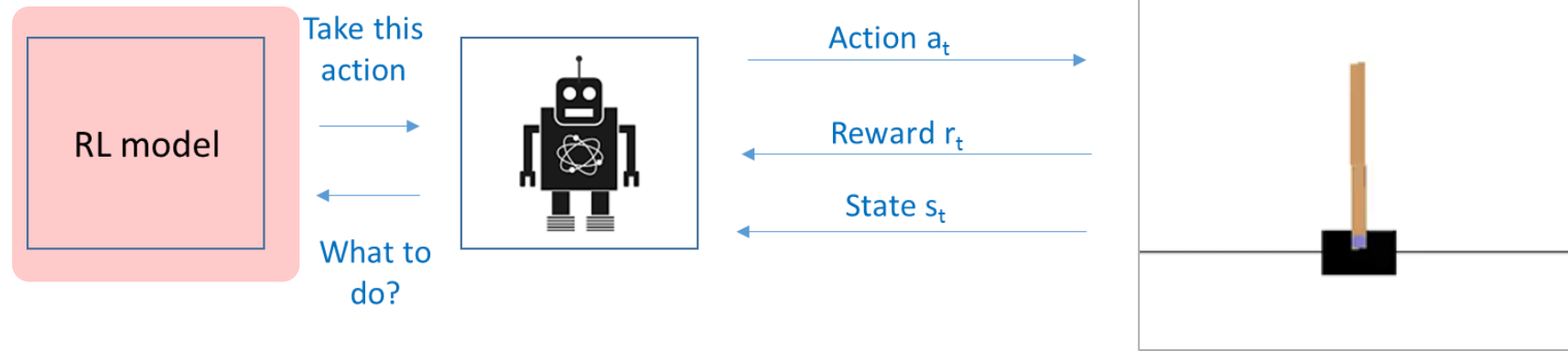
# gym returns values
next_state, reward, done, info = env.step(action)  # Asking to Gym box

# see the values
print("current state: ", current_state)
print("next state   : ", next_state)
print("reward       : ", reward)
print("finished?    : ", done)

current state: [ 0.02754178 -0.01649114 -0.01287597 -0.01059267]
next state   : [ 0.02721195  0.17881308 -0.01308783 -0.30731015]
reward       : 1.0
finished?    : False
```

- State: 4 dimension vector
  - [0]: cart position [-2.4 to 2.4]
  - [1]: cart velocity [-Inf to Inf]
  - [2]: pole angle [-42.8° to 41.8°]
  - [3]: pole velocity at tip [-Inf to Inf]

# Reinforcement Learning (RL) algorithm



- 1) Deep Q-Networks (DQN)
- 2) Policy Gradient (PG)
- 3) Actor Critic (AC)

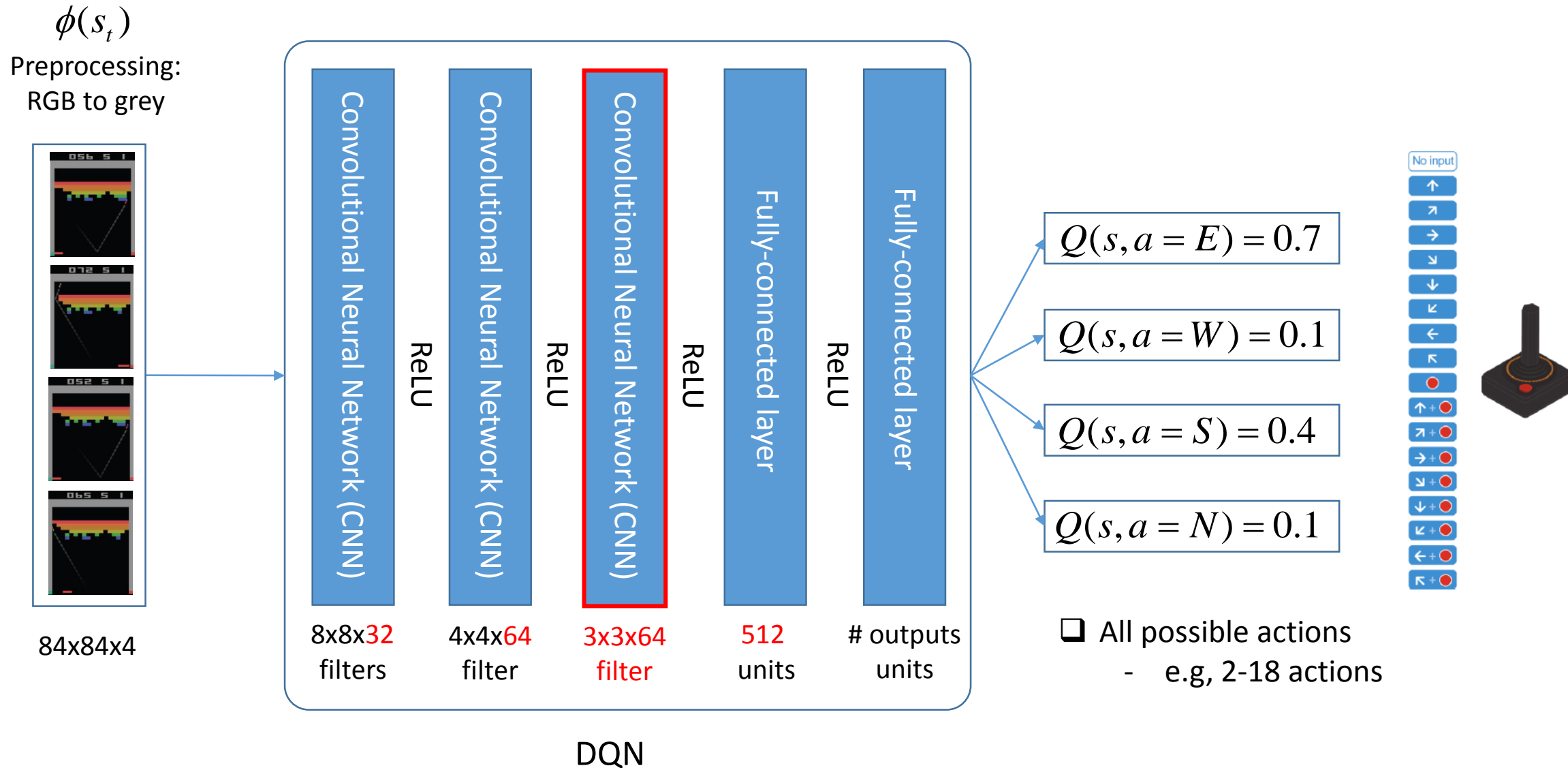


# DQN implementation

<https://github.com/hunkim/DeepLearningZeroToAll>

- 07\_3\_dqn\_2015\_cartpole.py

# DQN architecture (2015)



# Algorithm (DQN 2015)

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$  -----> Data pool size and initialization

Initialize action-value function  $Q$  with random weights  $\theta$  -----> Weight of 1<sup>st</sup> NN initialization

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$  -----> Weight of 2<sup>nd</sup> NN initialization

**For** episode = 1,  $M$  **do**

    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$  -----> Preprocessing, e.g., RGB to gray

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

        With probability  $\varepsilon$  select a random action  $a_t$   
        otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$  } -----> Action selection using E-greedy: off-policy

        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  $D$  -----> Experience replay

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

        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the -----> Update NN ( $\theta$ ) without changing NN ( $\theta^-$ )  
        network parameters  $\theta$

        Every  $C$  steps reset  $\hat{Q} = Q$  -----> Replace NN ( $\theta^-$ ) with NN ( $\theta$ ) every  $C$  steps

**End For**

**End For**

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

Initialize target action-value function  $\hat{Q}$

For episode = 1,  $M$  do

Initialize sequence  $s_1$

For  $t = 1, T$  do

With probability  $\epsilon$

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

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  $D$

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  $\theta$

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

End For

End For

### Data pool size and initialization

- deque(): list-like container with fast operation

- <https://docs.python.org/2/library/collections.html>

# store the previous observations in replay memory

replay\_buffer = deque(maxlen=REPLAY\_MEMORY)

last\_100\_game\_reward = deque(maxlen=100)

with tf.Session() as sess:

mainDQN = dqn.DQN(sess, INPUT\_SIZE, OUTPUT\_SIZE, name="main")

targetDQN = dqn.DQN(sess, INPUT\_SIZE, OUTPUT\_SIZE, name="target")

sess.run(tf.global\_variables\_initializer())

# initial copy q\_net -> target\_net

copy\_ops = get\_copy\_var\_ops(dest\_scope\_name="target",  
src\_scope\_name="main")

sess.run(copy\_ops)

for episode in range(MAX\_EPISODES):

e = 1. / ((episode / 10) + 1)

done = False

step\_count = 0

state = env.reset()

while not done:

if np.random.rand() < e:

action = env.action\_space.sample()

else:

# Choose an action by greedily from the Q-network

action = np.argmax(mainDQN.predict(state))

# Get new state and reward from environment

next\_state, reward, done, \_ = env.step(action)

if done: # Penalty

reward = -1

# Save the experience to our buffer

replay\_buffer.append((state, action, reward, next\_state, done))

if len(replay\_buffer) > BATCH\_SIZE:

minibatch = random.sample(replay\_buffer, BATCH\_SIZE)

loss, \_ = replay\_train(mainDQN, targetDQN, minibatch)

if step\_count % TARGET\_UPDATE\_FREQUENCY == 0:

sess.run(copy\_ops)

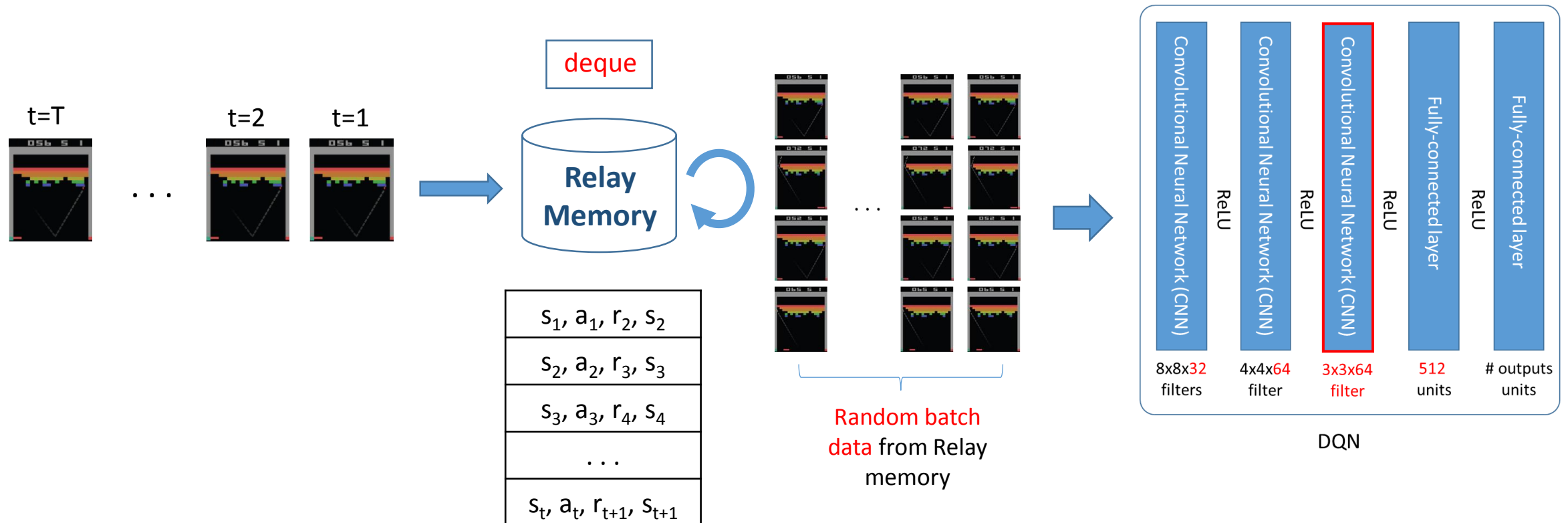
state = next\_state

step\_count += 1

print("Episode: {} steps: {}".format(episode, step\_count))

# DQN Code: experience replay

- ❑ Consecutive data frames are highly correlated
- ❑ Experience replay aims to remove the correlation between data samples



## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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

**For** episode = 1,  $M$  **do**

Initialize sequence  $s_1 = \{x_1\}$  and

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

With probability  $\epsilon$  select a random action  $a_t$

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

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  $D$

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  $\theta$

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

**End For**

**End For**

- Main Q-Net and target  $\hat{Q}$ -Net creation
- Copy from main Q-Net to target  $\hat{Q}$ -Net

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
mainDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="main")
targetDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="target")
sess.run(tf.global_variables_initializer())
```

```
# initial copy q_net -> target net
```

```
copy_ops = get_copy_var_ops(dest_scope_name="target",
                             src_scope_name="main")
sess.run(copy_ops)
```

```
for episode in range(MAX_EPISODES):
```

```
    e = 1. / ((episode / 10) + 1)
```

```
    done = False
```

```
    step_count = 0
```

```
    state = env.reset()
```

```
    while not done:
```

```
        if np.random.rand() < e:
```

```
            action = env.action_space.sample()
```

```
        else:
```

```
            # Choose an action by greedily from the Q-network
```

```
            action = np.argmax(mainDQN.predict(state))
```

```
    # Get new state and reward from environment
```

```
    next_state, reward, done, _ = env.step(action)
```

```
    if done: # Penalty
```

```
        reward = -1
```

```
    # Save the experience to our buffer
```

```
    replay_buffer.append((state, action, reward, next_state, done))
```

```
    if len(replay_buffer) > BATCH_SIZE:
```

```
        minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
        loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
    if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

```
        sess.run(copy_ops)
```

```
    state = next_state
```

```
    step_count += 1
```

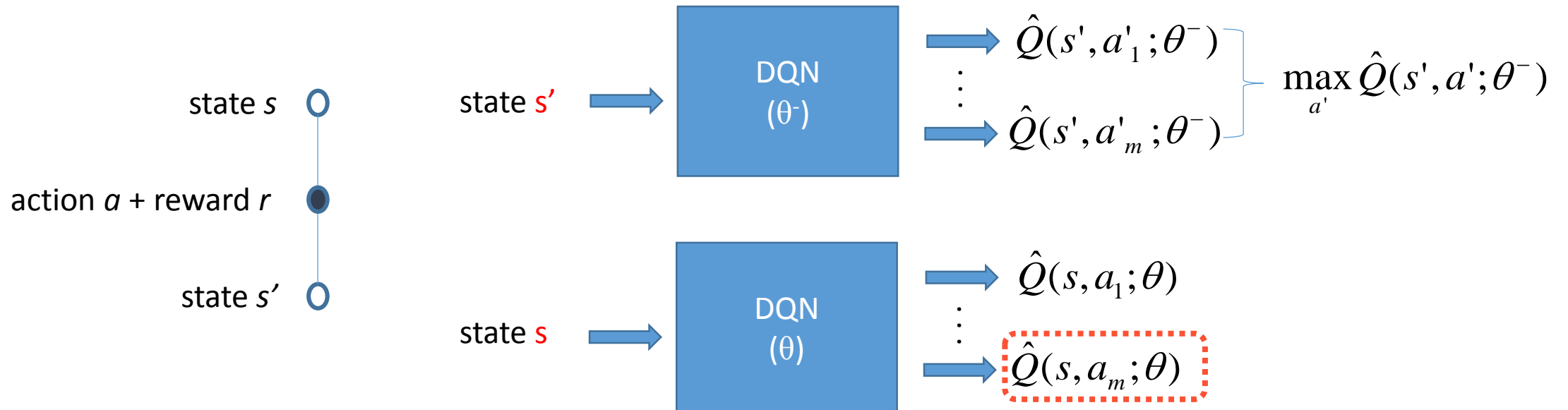
```
print("Episode: {} steps: {}".format(episode, step_count))
```



# DQN Code: fixed Q-target

$$L_i(\theta_i) = E_{(s,a,r,s') \sim U(D)} \left[ \left( \underbrace{r + \gamma \max_{a'} Q(s', a' | \theta_i^-)}_{\text{Target Q network}} - \underbrace{Q(s, a, \theta_i)}_{\text{Main Q network}} \right)^2 \right]$$

- Sample data  $(s,a,r,s')$  randomly drawn from data pool  $U(D)$
- Experience Replay
- Two different neural networks
- Fixed Q-target



## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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

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

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  $\theta$

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

**End For**

**End For**

▪ Decaying E-greedy implementation

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
    mainDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="main")
```

```
    targetDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="target")
```

```
    sess.run(tf.global_variables_initializer())
```

```
# initial copy q_net -> target_net
```

```
copy_ops = get_copy_var_ops(dest_scope_name="target",
                             src_scope_name="main")
```

```
sess.run(copy_ops)
```

```
for episode in range(MAX_EPISODES):
```

```
    e = 1. / ((episode / 10) + 1)
```

```
    done = False
```

```
    step_count = 0
```

```
    state = env.reset()
```

```
    while not done:
```

```
        if np.random.rand() < e:
            action = env.action_space.sample()
```

```
        else:
```

```
            # Choose an action by greedily from the Q-network
            action = np.argmax(mainDQN.predict(state))
```

```
# Get new state and reward from environment
```

```
next_state, reward, done, _ = env.step(action)
```

```
if done: # Penalty
```

```
    reward = -1
```

```
# Save the experience to our buffer
```

```
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
```

```
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

```
    sess.run(copy_ops)
```

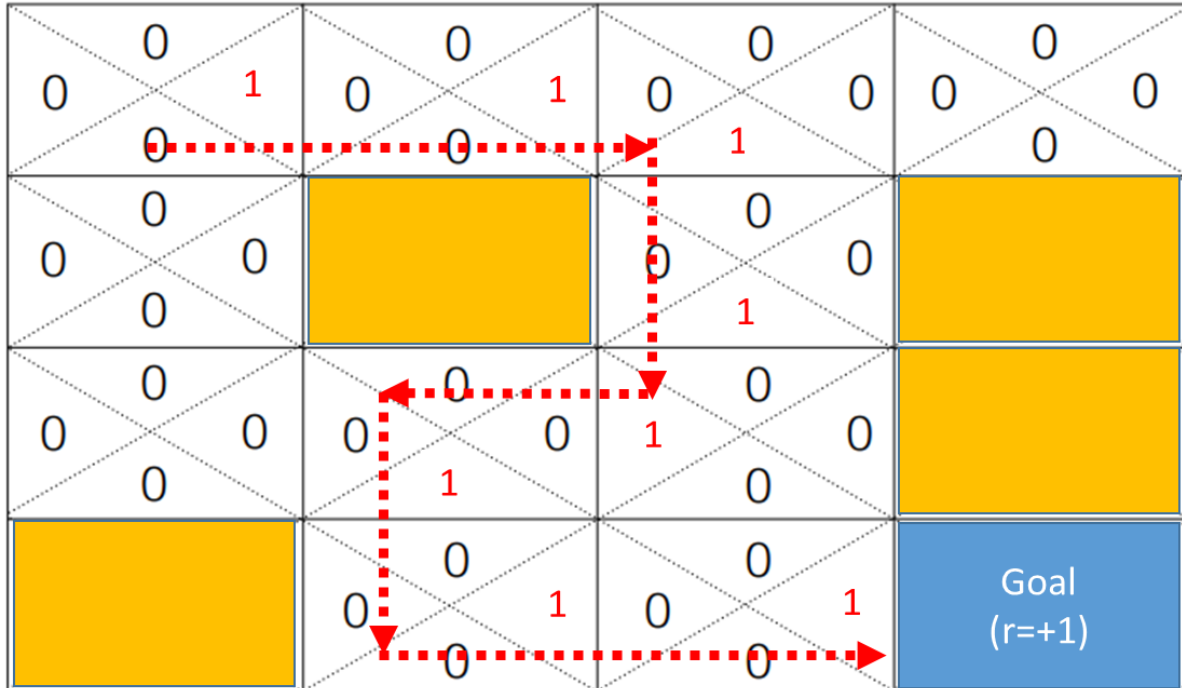
```
state = next_state
```

```
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```



# DQN Code: Exploit vs Exploration



## Decaying E-greedy policy

```
for i in range (1000)
```

$e = 0.1 / (i+1)$  → Random to deterministic decision as iteration goes on

```
if rand < e:
```

```
action = random
```

```
else:
```

$$\text{action} = \text{argmax}(Q(s,a))$$

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  select a random action  $a_t$

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

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  $D$

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

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  $\theta$

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

**End For**

**End For**

■ OpenAI GYM emulator

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
    mainDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="main")
```

```
    targetDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="target")
```

```
    sess.run(tf.global_variables_initializer())
```

```
# initial copy q_net -> target_net
```

```
copy_ops = get_copy_var_ops(dest_scope_name="target",
                             src_scope_name="main")
```

```
sess.run(copy_ops)
```

```
for episode in range(MAX_EPISODES):
```

```
    e = 1. / ((episode / 10) + 1)
```

```
    done = False
```

```
    step_count = 0
```

```
    state = env.reset()
```

```
    while not done:
```

```
        if np.random.rand() < e:
```

```
            action = env.action_space.sample()
```

```
        else:
```

```
            # Choose an action by greedily from the Q-network
```

```
            action = np.argmax(mainDQN.predict(state))
```

```
# Get new state and reward from environment
```

```
next_state, reward, done, _ = env.step(action)
```

```
if done: # Penalty
```

```
    reward = -1
```

```
# Save the experience to our buffer
```

```
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
```

```
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

```
    sess.run(copy_ops)
```

```
state = next_state
```

```
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```

## ❑ Action:

- 0: left
- 1: right

## ❑ State: 4 dimension vector

- [0]: cart position [-2.4 to 2.4]
- [1]: cart velocity [-Inf to Inf]
- [2]: pole angle [-42.8° to 41.8°]
- [3]: pole velocity at tip [-Inf to Inf]

## ❑ Initial state

- Random values between  $\pm 0.05$

## ❑ Reward:

- +1 each unit time if it is not fallen

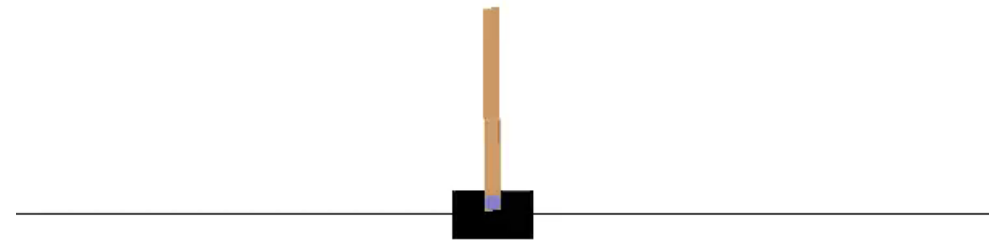
## ❑ Episode Termination

- Pole Angle is more than  $\pm 12^\circ$
- Cart Position is more than  $\pm 2.4$
- Episode length is greater than 200 (unit time)

## ❑ Training Termination

- Average reward is greater than or equal to 195 over 100 consecutive trials

```
# Get new state and reward from environment
next_state, reward, done, _ = env.step(action)
```



```
10 = {ndarray} [-0.1062376 -0.37924474 0.19689548 0.9208114 ]
```

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  select a random action  $a_t$

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

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  $D$

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  $\theta$

- Experience replay implementation
- Wait until collecting enough number of input data (batch\_size)
- Then, randomly sample "batch\_size" number of inputs from the buffer

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
    mainDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="main")
```

```
    targetDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="target")
```

```
    sess.run(tf.global_variables_initializer())
```

```
# initial copy q_net -> target_net
```

```
copy_ops = get_copy_var_ops(dest_scope_name="target",
                             src_scope_name="main")
```

```
sess.run(copy_ops)
```

```
for episode in range(MAX_EPISODES):
```

```
    e = 1. / ((episode / 10) + 1)
```

```
    done = False
```

```
    step_count = 0
```

```
    state = env.reset()
```

```
    while not done:
```

```
        if np.random.rand() < e:
```

```
            action = env.action_space.sample()
```

```
        else:
```

```
            # Choose an action by greedily from the Q-network
```

```
            action = np.argmax(mainDQN.predict(state))
```

```
# Get new state and reward from environment
```

```
next_state, reward, done, _ = env.step(action)
```

```
if done: # Penalty
```

```
    reward = -1
```

```
# Save the experience to our buffer
```

```
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
```

```
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

```
    sess.run(copy_ops)
```

```
state = next_state
```

```
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```



## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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  $\theta$

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

**End For**

**End For**

■ Training the main Q-Net given minibatch input

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
    mainDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="main")
```

```
    targetDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="target")
```

```
    sess.run(tf.global_variables_initializer())
```

```
# initial copy q_net -> target_net
```

```
copy_ops = get_copy_var_ops(dest_scope_name="target",
                             src_scope_name="main")
```

```
sess.run(copy_ops)
```

```
for episode in range(MAX_EPISODES):
```

```
    e = 1. / ((episode / 10) + 1)
```

```
    done = False
```

```
    step_count = 0
```

```
    state = env.reset()
```

```
    while not done:
```

```
        if np.random.rand() < e:
```

```
            action = env.action_space.sample()
```

```
        else:
```

```
            # Choose an action by greedily from the Q-network
```

```
            action = np.argmax(mainDQN.predict(state))
```

```
# Get new state and reward from environment
```

```
next_state, reward, done, _ = env.step(action)
```

```
if done: # Penalty
```

```
    reward = -1
```

```
# Save the experience to our buffer
```

```
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
```

```
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

```
    sess.run(copy_ops)
```

```
state = next_state
```

```
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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  $\theta$

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

**End For**

**End For**

- Training the main Q-Net given minibatch input

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
def replay_train(mainDQN: dqn.DQN, targetDQN: dqn.DQN, train_batch: list) -> float:
    """Trains 'mainDQN' with target Q values given by 'targetDQN'"""
```

```
Args:
```

```
mainDQN (dqn.DQN): Main DQN that will be trained
targetDQN (dqn.DQN): Target DQN that will predict Q_target
train_batch (list): Minibatch of replay memory
    Each element is (s, a, r, s', done)
    [(state, action, reward, next_state, done), ...]
```

```
Returns:
```

```
float: After updating 'mainDQN', it returns a 'loss'"""
```

```
states = np.vstack([x[0] for x in train_batch])
actions = np.array([x[1] for x in train_batch])
rewards = np.array([x[2] for x in train_batch])
next_states = np.vstack([x[3] for x in train_batch])
done = np.array([x[4] for x in train_batch])
```

```
X = states
```

```
Q_target = rewards + DISCOUNT_RATE * np.max(targetDQN.predict(next_states), axis=1) * ~done
```

```
y = mainDQN.predict(states)
y[np.arange(len(X)), actions] = Q_target
```

```
# Train our network using target and predicted Q values on each episode
```

```
return mainDQN.update(X, y)
```

```
reward = -1
```

```
# Save the experience to our buffer
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
    sess.run(copy_ops)
```

```
state = next_state
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```



## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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  $\theta$

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

**End For**

**End For**

- Training the main Q-Net given minibatch input

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
def replay_train(mainDQN: dqn.DQN, targetDQN: dqn.DQN, train_batch: list) -> float:
    """Trains 'mainDQN' with target Q values given by 'targetDQN'"""
```

```
Args:
```

```
mainDQN (dqn.DQN): Main DQN that will be trained
targetDQN (dqn.DQN): Target DQN that will predict Q_target
train_batch (list): Minibatch of replay memory
    Each element is (s, a, r, s', done)
    [(state, action, reward, next_state, done), ...]
```

```
Returns:
```

```
float: After updating 'mainDQN', it returns a 'loss'"""
```

```
states = np.vstack([x[0] for x in train_batch])
actions = np.array([x[1] for x in train_batch])
rewards = np.array([x[2] for x in train_batch])
next_states = np.vstack([x[3] for x in train_batch])
done = np.array([x[4] for x in train_batch])
```

```
X = states
```

```
Q_target = rewards + DISCOUNT_RATE * np.max(targetDQN.predict(next_states), axis=1) * ~done
```

```
y = mainDQN.predict(states)
y[np.arange(len(X)), actions] = Q_target
```

```
# Train our network using target and predicted Q values on each episode
return mainDQN.update(X, y)
```

```
reward = -1
```

```
# Save the experience to our buffer
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
    sess.run(copy_ops)
```

```
state = next_state
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```

## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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  $\theta$

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

**End For**

**End For**

■ Training the main Q-Net given minibatch input

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
def replay_train(mainDQN: dqn.DQN, targetDQN: dqn.DQN, train_batch: list) -> float:
    """Trains 'mainDQN' with target Q values given by 'targetDQN'"""
```

```
Args:
```

```
mainDQN (dqn.DQN): Main DQN that will be trained
targetDQN (dqn.DQN): Target DQN that will predict Q_target
train_batch (list): Minibatch of replay memory
    Each element is (s, a, r, s', done)
    [(state, action, reward, next_state, done), ...]
```

```
Returns:
```

```
float: After updating 'mainDQN', it returns a 'loss'"""
```

```
states = np.vstack([x[0] for x in train_batch])
actions = np.array([x[1] for x in train_batch])
rewards = np.array([x[2] for x in train_batch])
next_states = np.vstack([x[3] for x in train_batch])
done = np.array([x[4] for x in train_batch])
```

```
X = states
```

```
Q_target = rewards + DISCOUNT_RATE * np.max(targetDQN.predict(next_states), axis=1) * ~done
```

```
y = mainDQN.predict(states)
y[np.arange(len(X)), actions] = Q_target
```

```
# Train our network using target and predicted Q values on each episode
```

```
return mainDQN.update(X, y)
```

```
reward = -1
```

```
# Save the experience to our buffer
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
    sess.run(copy_ops)
```

```
state = next_state
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```



## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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  $\theta$

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

**End For**

**End For**

▪ Copy NN ( $\theta^-$ ) with NN ( $\theta$ ) every  $C$  steps

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
    mainDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="main")
```

```
    targetDQN = dqn.DQN(sess, INPUT_SIZE, OUTPUT_SIZE, name="target")
```

```
    sess.run(tf.global_variables_initializer())
```

```
# initial copy q_net -> target_net
```

```
copy_ops = get_copy_var_ops(dest_scope_name="target",
                             src_scope_name="main")
```

```
sess.run(copy_ops)
```

```
for episode in range(MAX_EPISODES):
```

```
    e = 1. / ((episode / 10) + 1)
```

```
    done = False
```

```
    step_count = 0
```

```
    state = env.reset()
```

```
    while not done:
```

```
        if np.random.rand() < e:
```

```
            action = env.action_space.sample()
```

```
        else:
```

```
            # Choose an action by greedily from the Q-network
```

```
            action = np.argmax(mainDQN.predict(state))
```

```
# Get new state and reward from environment
```

```
next_state, reward, done, _ = env.step(action)
```

```
if done: # Penalty
```

```
    reward = -1
```

```
# Save the experience to our buffer
```

```
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
```

```
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

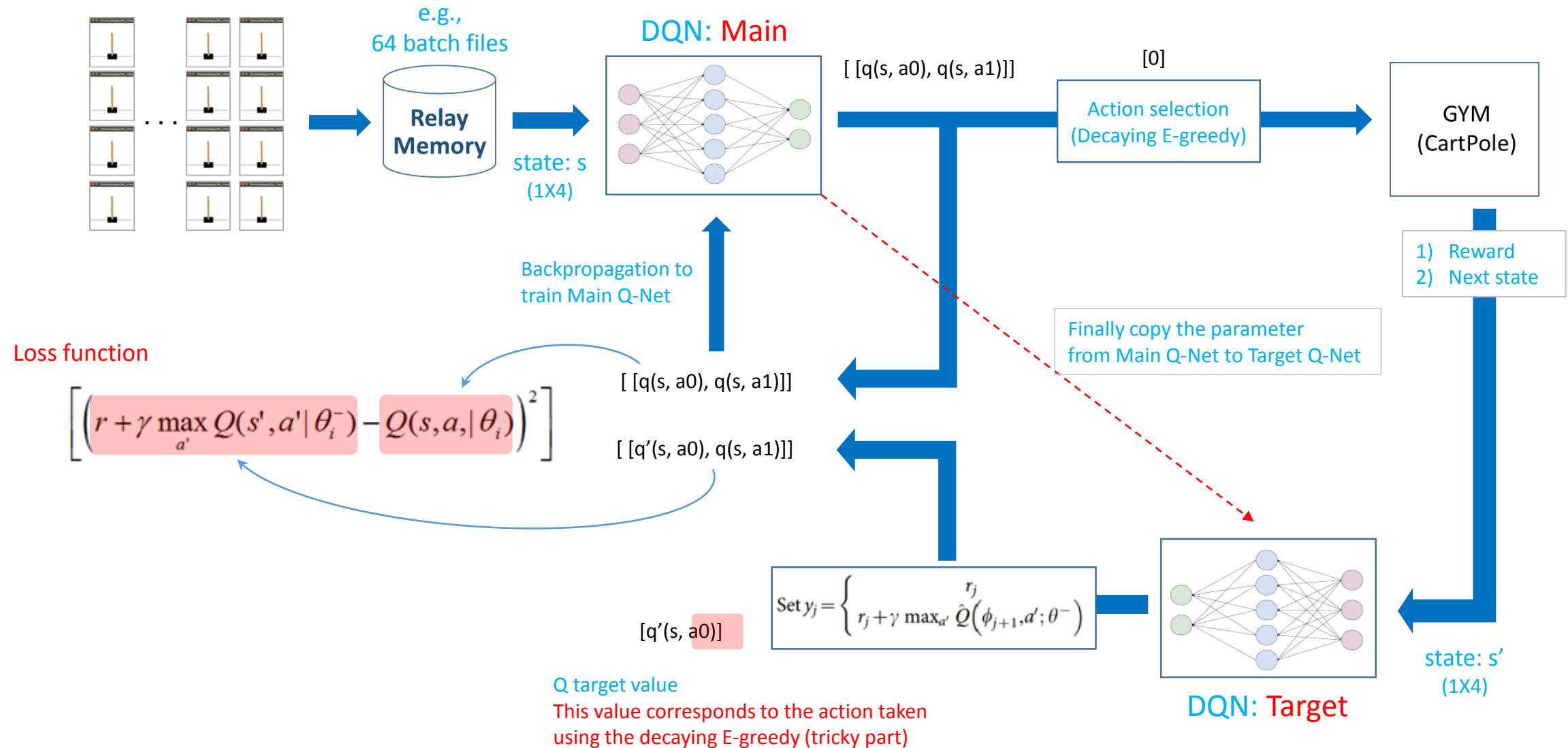
```
    sess.run(copy_ops)
```

```
    state = next_state
```

```
    step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```

# Big picture for the implementation of DQN



# PG implementation

<http://karpathy.github.io/2016/05/31/rl/>

<https://github.com/hunkim/DeepLearningZeroToAll>

- 08\_1\_pg\_cartpole.py

---

**Algorithm 1** “Vanilla” policy gradient algorithm
 

---

Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration=1, 2, ... **do**

Collect a set of trajectories by executing the current policy

At each timestep in each trajectory, compute

the *return*  $R_t = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$ , and

the *advantage estimate*  $\hat{A}_t = R_t - b(s_t)$ .

Re-fit the baseline, by minimizing  $\|b(s_t) - R_t\|^2$ ,

summed over all trajectories and timesteps.

Update the policy, using a policy gradient estimate  $\hat{g}$ ,

which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$

**end for**

---

```
for step in range(max_num_episodes):
    # Initialize x stack, y stack, and rewards
    xs = np.empty(shape=[0, input_size])
    ys = np.empty(shape=[0, 1])
    rewards = np.empty(shape=[0, 1])

    reward_sum = 0
    observation = env.reset()

    while True:
        x = np.reshape(observation, [1, input_size])

        # Run the neural net to determine output
        action_prob = sess.run(action_pred, feed_dict={X: x})

        # Determine the output based on our net, allowing for some randomness
        action = 0 if action_prob < np.random.uniform() else 1

        # Append the observations and outputs for learning
        xs = np.vstack([xs, x])
        ys = np.vstack([ys, action]) # Fake action

        # Determine the outcome of our action
        observation, reward, done, _ = env.step(action)
        rewards = np.vstack([rewards, reward])
        reward_sum += reward

        if done:
            # Determine standardized rewards
            discounted_rewards = discount_rewards(rewards)
            # Normalization
            discounted_rewards = (discounted_rewards - discounted_rewards.mean())
                                / (discounted_rewards.std() + 1e-7)
            l, _ = sess.run([loss, train],
                           feed_dict={X: xs, Y: ys, advantages: discounted_rewards})
```



## Algorithm 1 “Vanilla” policy gradient algorithm

Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration=1, 2, ... **do**

Collect a set of trajectories by executing the current policy

At each timestep in each trajectory, compute

the return  $R_t = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$ , and

the advantage estimate  $\hat{A}_t = R_t - b(s_t)$ .

Re-fit the baseline, by minimizing  $\|b(s_t) - R_t\|^2$ ,

summed over all tra

Update the policy, us

which is a sum of t

**end for**

“np.vstack” stacks each event on a trajectory until an episode finish

```
for step in range(max_num_episodes):
    # Initialize x stack, y stack, and rewards
    xs = np.empty(shape=[0, input_size])
    ys = np.empty(shape=[0, 1])
    rewards = np.empty(shape=[0, 1])

    reward_sum = 0
    observation = env.reset()

    while True:
        x = np.reshape(observation, [1, input_size])

        # Run the neural net to determine output
        action_prob = sess.run(action_pred, feed_dict={X: x})

        # Determine the output based on our net, allowing for some randomness
        action = 0 if action_prob < np.random.uniform() else 1

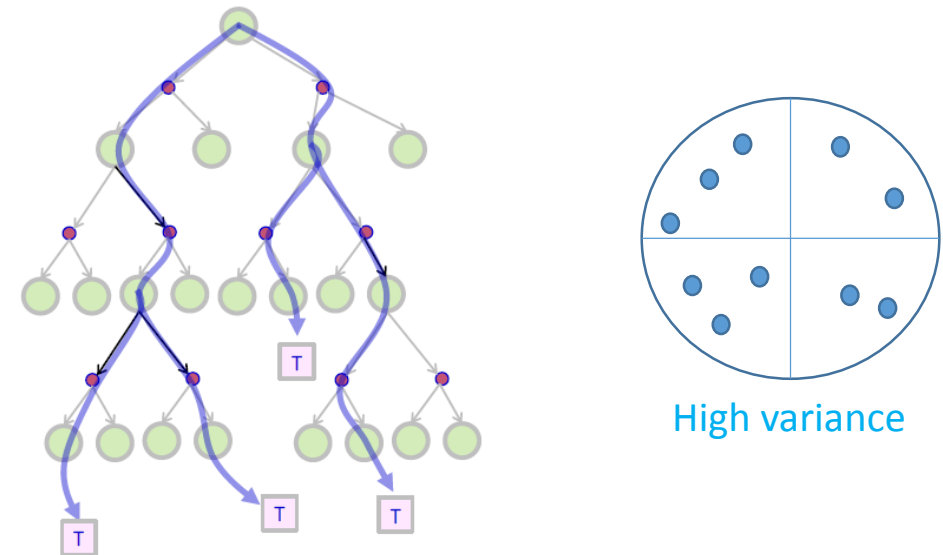
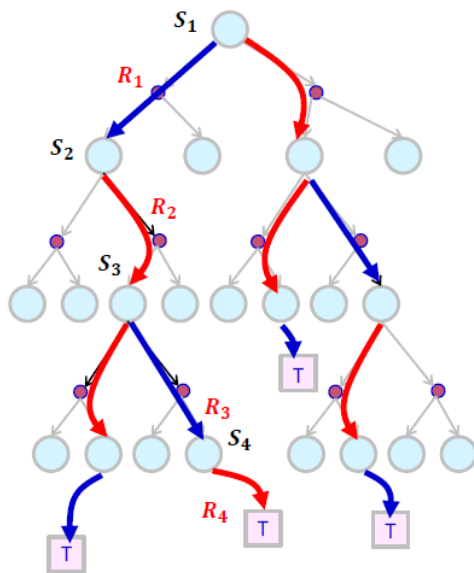
        # Append the observations and outputs for learning
        xs = np.vstack([xs, x])
        ys = np.vstack([ys, action]) # Fake action

        # Determine the outcome of our action
        observation, reward, done, _ = env.step(action)
        rewards = np.vstack([rewards, reward])
        reward_sum += reward

        if done:
            # Determine standardized rewards
            discounted_rewards = discount_rewards(rewards)
            # Normalization
            discounted_rewards = (discounted_rewards - discounted_rewards.mean())
                                / (discounted_rewards.std() + 1e-7)
            l, _ = sess.run([loss, train],
                           feed_dict={X: xs, Y: ys, advantages: discounted_rewards})
```

# PG Code: Q learning vs Policy gradient

Q learning	Policy gradient
<ul style="list-style-type: none"><li>■ Learning <math>Q(s,a)</math>: modeling (Reward) values of actions<ul style="list-style-type: none"><li>- Value based approach: <b>learning Q values</b></li></ul></li></ul>	<ul style="list-style-type: none"><li>■ Learning <math>\pi(a)</math>: modeling probability of actions<ul style="list-style-type: none"><li>- Policy based approach: <b>learning policy directly</b></li></ul></li></ul>
<ul style="list-style-type: none"><li>■ Deterministic policies:<ul style="list-style-type: none"><li>- e.g., <b>cannot</b> model rock-paper-scissors game</li></ul></li></ul>	<ul style="list-style-type: none"><li>■ Stochastic policies<ul style="list-style-type: none"><li>- e.g., <b>can</b> model rock-paper-scissors game</li></ul></li></ul>
<ul style="list-style-type: none"><li>■ <b>Off-policy</b>: an action is taken greedily<ul style="list-style-type: none"><li>- Greedy search to calculate <math>Q(s,a)</math> and then determine a policy</li></ul></li></ul>	<ul style="list-style-type: none"><li>■ <b>On-policy</b>: an action is taken with a policy<ul style="list-style-type: none"><li>- Following a trajectory created by a policy and update it with given reward at the end.</li></ul></li></ul>
<ul style="list-style-type: none"><li>■ Learning update occurred step-by-step (bootstrapping)<ul style="list-style-type: none"><li>- Low variance but high bias</li></ul></li></ul>	<ul style="list-style-type: none"><li>■ Learning update occurred episode-by-episode<ul style="list-style-type: none"><li>- High variance but low bias</li></ul></li></ul>



## Algorithm 1 “Vanilla” policy gradient algorithm

Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration=1, 2, ... **do**

Collect a set of trajectories by executing the current policy

At each timestep in each trajectory, compute the return  $R_t = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$ , and the advantage estimate  $\hat{A}_t = R_t - b(s_t)$ .

Re-fit the baseline, by minimizing  $\|b(s_t) - R_t\|^2$ , summed over all trajectories and timesteps.

Update the policy, using a policy gradient estimate  $\hat{g}$ , which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$

**end for**

If an episode finishes,

- 1) Discounted reward is calculated
- 2) Then, it is normalized
  - <http://karpathy.github.io/2016/05/31/rl/>

```
for step in range(max_num_episodes):
    # Initialize x stack, y stack, and rewards
    xs = np.empty(shape=[0, input_size])
    ys = np.empty(shape=[0, 1])
    rewards = np.empty(shape=[0, 1])

    reward_sum = 0
    observation = env.reset()

    while True:
        x = np.reshape(observation, [1, input_size])

        # Run the neural net to determine output
        action_prob = sess.run(action_pred, feed_dict={X: x})

        # Determine the output based on our net, allowing for some randomness
        action = 0 if action_prob < np.random.uniform() else 1

        # Append the observations and outputs for learning
        xs = np.vstack([xs, x])
        ys = np.vstack([ys, action]) # Fake action

        # Determine the outcome of our action
        observation, reward, done, _ = env.step(action)
        rewards = np.vstack([rewards, reward])
        reward_sum += reward

        if done:
            # Determine standardized rewards
            discounted_rewards = discount_rewards(rewards)
            # Normalization
            discounted_rewards = (discounted_rewards - discounted_rewards.mean())
                                / (discounted_rewards.std() + 1e-7)
            l, _ = sess.run([loss, train],
                           feed_dict={X: xs, Y: ys, advantages: discounted_rewards})
```



## Algorithm 1 “Vanilla” policy gradient algorithm

Initialize policy parameter  $\theta$ , baseline  $b$

**for** iteration=1, 2, ... **do**

Collect a set of trajectories by executing the current policy

At each timestep in each trajectory, compute

the return  $R_t = \sum_{t'=t}^{T-1} \gamma^{t'-t} r_{t'}$ , and

the advantage estimate  $\hat{A}_t = R_t - b(s_t)$ .

Re-fit the baseline, by minimizing  $\|b(s_t) - R_t\|^2$ ,  
summed over all trajectories and timesteps.

Update the policy, using a policy gradient estimate  $\hat{g}$ ,  
which is a sum of terms  $\nabla_{\theta} \log \pi(a_t | s_t, \theta) \hat{A}_t$

**end for**

```
# Loss function: log_likelihood * advantages
log_lik = -Y*tf.log(action_pred) - (1 - Y)*tf.log(1 - action_pred)
loss = tf.reduce_sum(log_lik * advantages)

# Learning
train = tf.train.AdamOptimizer(learning_rate=learning_rate).minimize(loss)
```

```
for step in range(max_num_episodes):
    # Initialize x stack, y stack, and rewards
    xs = np.empty(shape=[0, input_size])
    ys = np.empty(shape=[0, 1])
    rewards = np.empty(shape=[0, 1])

    reward_sum = 0
    observation = env.reset()

    while True:
        x = np.reshape(observation, [1, input_size])

        # Run the neural net to determine output
        action_prob = sess.run(action_pred, feed_dict={X: x})

        # Determine the output based on our net, allowing for some randomness
        action = 0 if action_prob < np.random.uniform() else 1

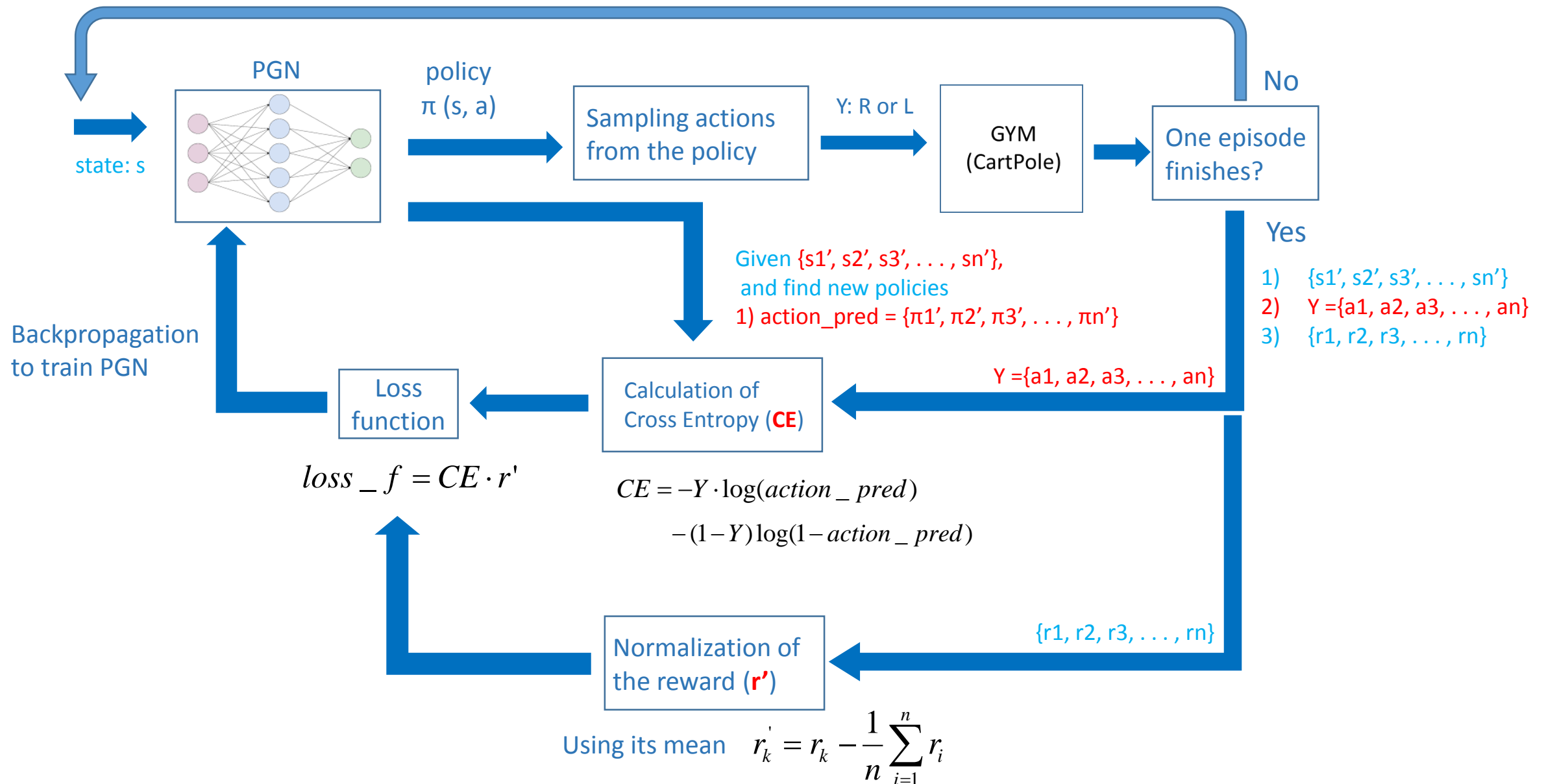
        # Append the observations and outputs for learning
        xs = np.vstack([xs, x])
        ys = np.vstack([ys, action]) # Fake action

        # Determine the outcome of our action
        observation, reward, done, _ = env.step(action)
        rewards = np.vstack([rewards, reward])
        reward_sum += reward

        if done:
            # Determine standardized rewards
            discounted_rewards = discount_rewards(rewards)
            # Normalization
            discounted_rewards = (discounted_rewards - discounted_rewards.mean())
                                / (discounted_rewards.std() + 1e-7)
            l, _ = sess.run([loss, train],
                           feed_dict={X: xs, Y: ys, advantages: discounted_rewards})
```



# Big picture for the implementation of PG

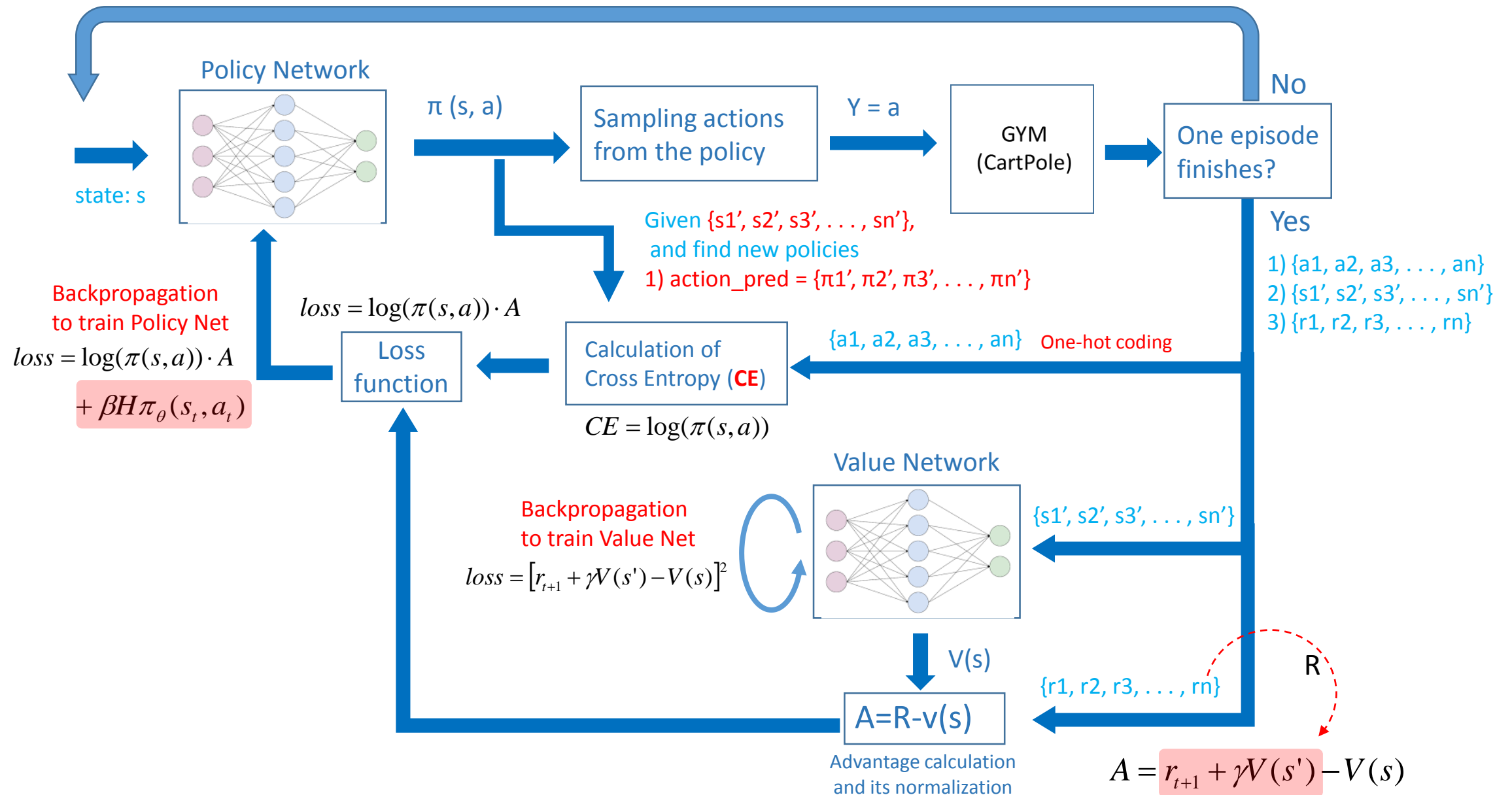


# Actor Advantage Critic (A2C)

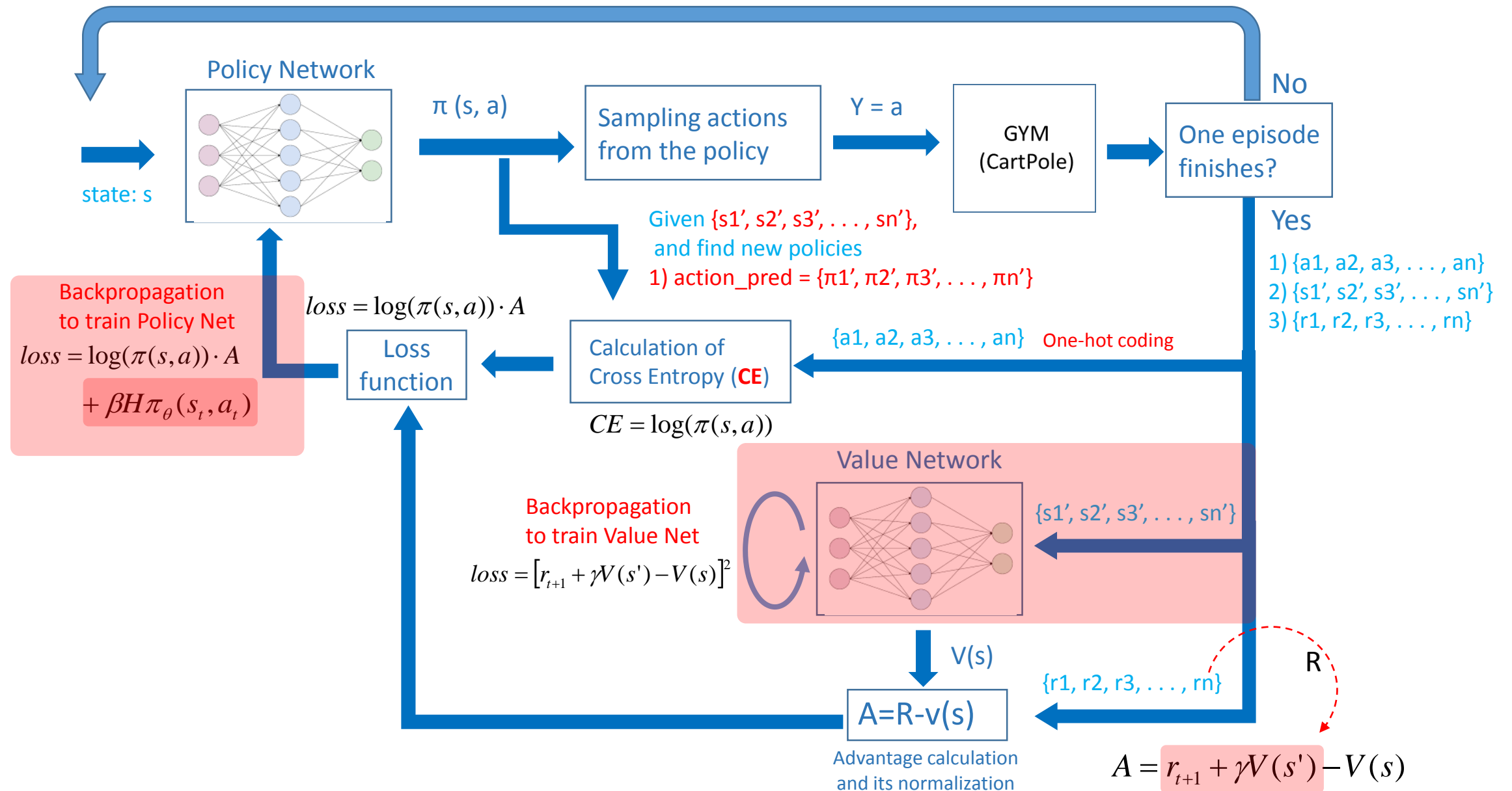
<https://github.com/hunkim/DeepLearningZeroToAll>

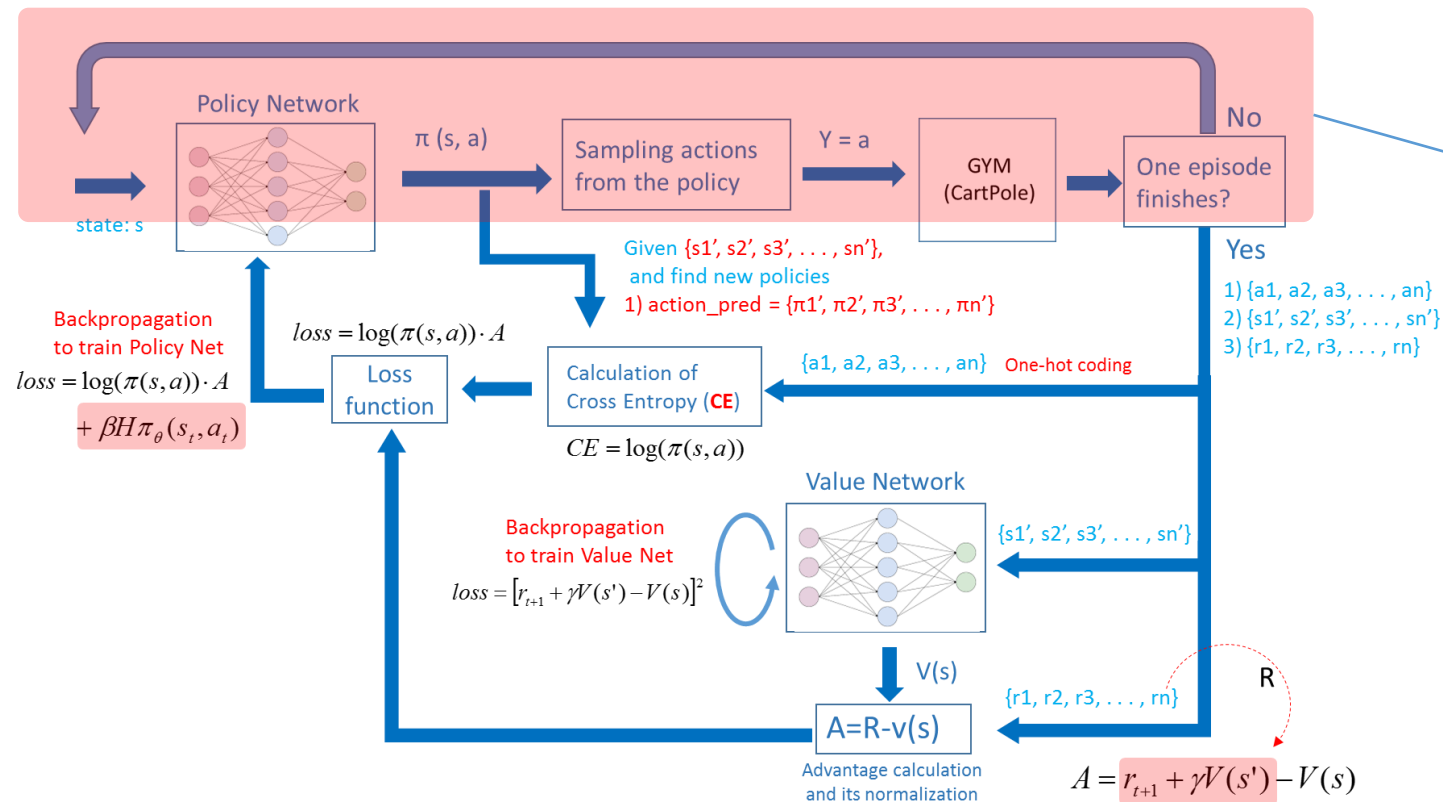
- 10\_1\_Actor\_Critic.ipynb

# Big picture for the implementation of AC



# Big picture for the implementation of AC





```
for episode in range(MAX_EPISODES):
    s = env.reset()
    done = False
```

```
s_list = []
a_list = []
r_list = []
```

```
episode_r = 0
```

Until one episode finishes

```
while not done:
```

```
    s = preprocess_state(s)
    a = agent.choose_an_action(s)

    s2, r, done, info = env.step(a)

    s_list.append(s)
    a_list.append(a)
    r_list.append(r)

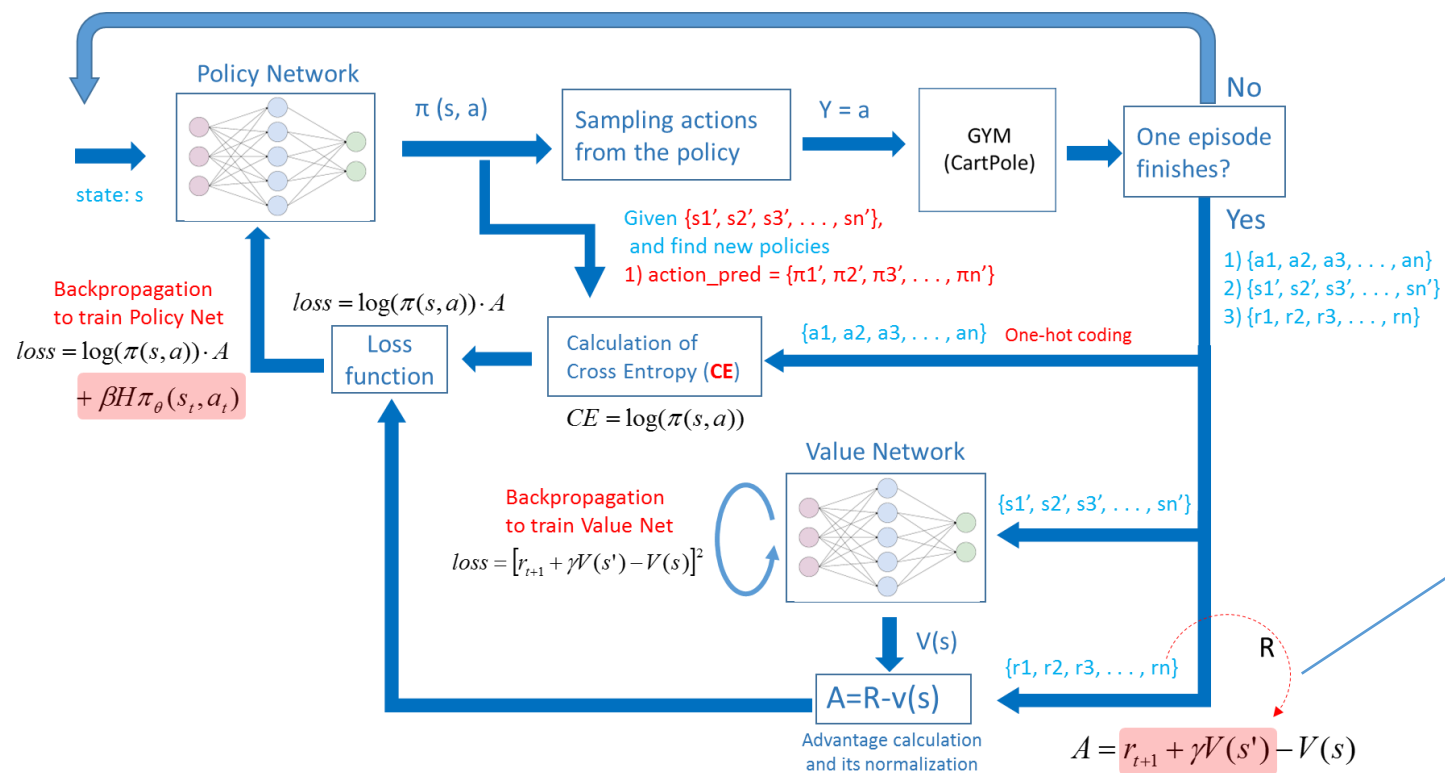
    s = s2

    episode_r += r
```

```
    a_list = preprocess_action(a_list, action_n)
```

```
    agent.train(np.vstack(s_list), a_list, r_list)
```

One-hot encoding for cross entropy calculation



```
def train(self, S, A, R):
    """ Train the actor critic networks

    (1) Compute discounted rewards R
    (2) Compute advantage values A = R - V
    (3) Perform gradients updates

    """

    def discount_rewards(r, gamma=0.99):
        """ take 1D float array of rewards and compute discounted rewards """
        discounted_r = np.zeros_like(r, dtype=np.float32)
        running_add = 0
        for t in reversed(range(len(r))):
            running_add = running_add * gamma + r[t]
            discounted_r[t] = running_add

        return discounted_r

    # 1. Get discounted 'R's
    R = discount_rewards(R)

    # 2. Get values
    V = self.value_net(S)

    # 3. Compute advantage
    A = R - V

    # 4. Compute loss
    loss = self.policy_net.loss(S, A)

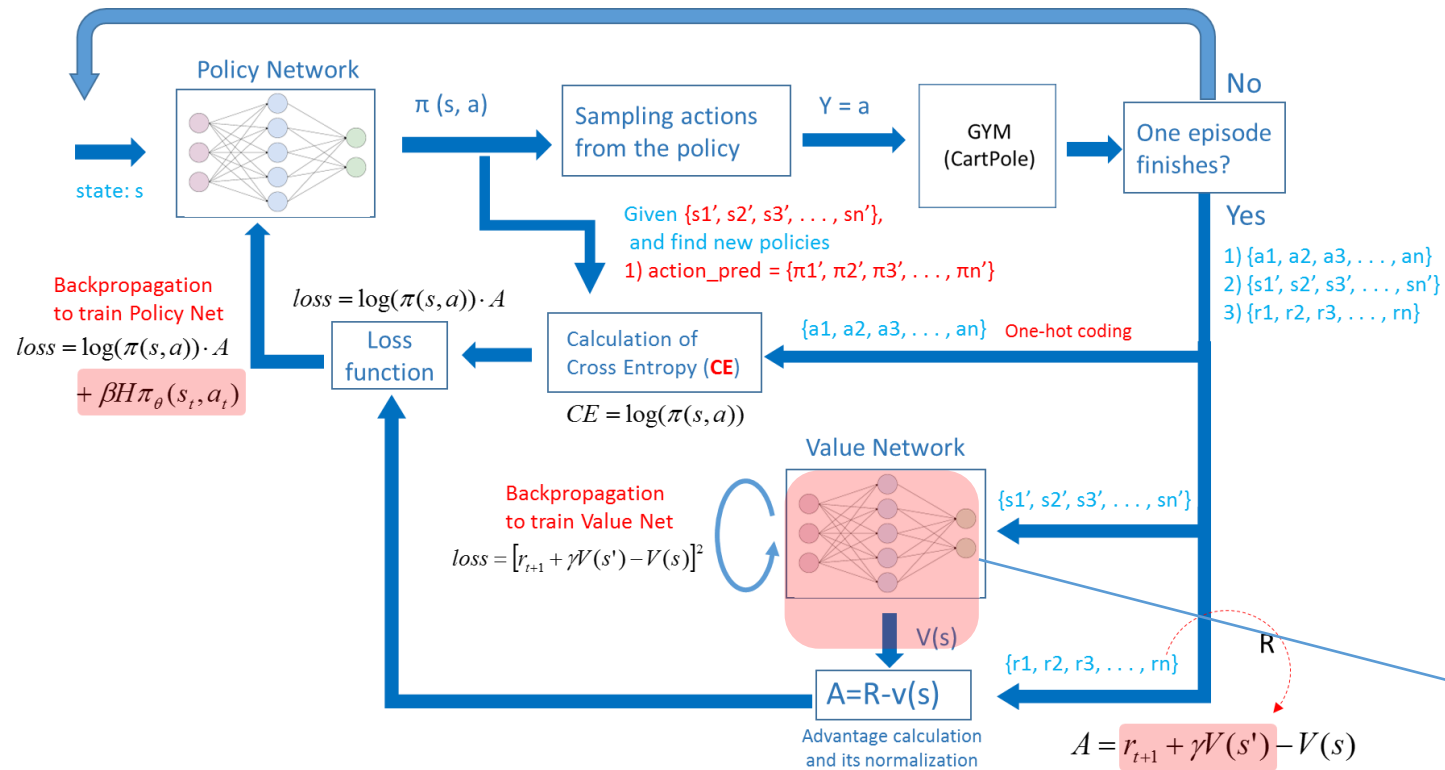
    # 5. Backpropagate
    self.optimizer.zero_grad()
    loss.backward()
    self.optimizer.step()

    # 6. Update policy and value networks
    self.policy_net.update()
    self.value_net.update()

    return loss
```

$R = [v(s_1'), v(s_2'), \dots, v(s_n')]$   
It represents value of each state then we can obtain  $(r + \gamma v(s'))$   
e.g.,

$s_1'$	$v(s_1')$	{float32} 10.466174
$s_2'$	$v(s_2')$	:
$s_3'$	$v(s_3')$	:
		{float32} 2.9701
		{float32} 1.99
$s_n'$	$v(s_n')$	{float32} 1.0



```
def train(self, S, A, R):
    """ Train the actor critic networks

    (1) Compute discounted rewards R
    (2) Compute advantage values A = R - V
    (3) Perform gradients updates

    """

    def discount_rewards(r, gamma=0.99):
        """ take 1D float array of rewards and compute discounted rewards """
        discounted_r = np.zeros_like(r, dtype=np.float32)
        running_add = 0

        for t in reversed(range(len(r))):
            running_add = running_add * gamma + r[t]
            discounted_r[t] = running_add

        return discounted_r

    # 1. Get discounted 'R's
    R = discount_rewards(R)

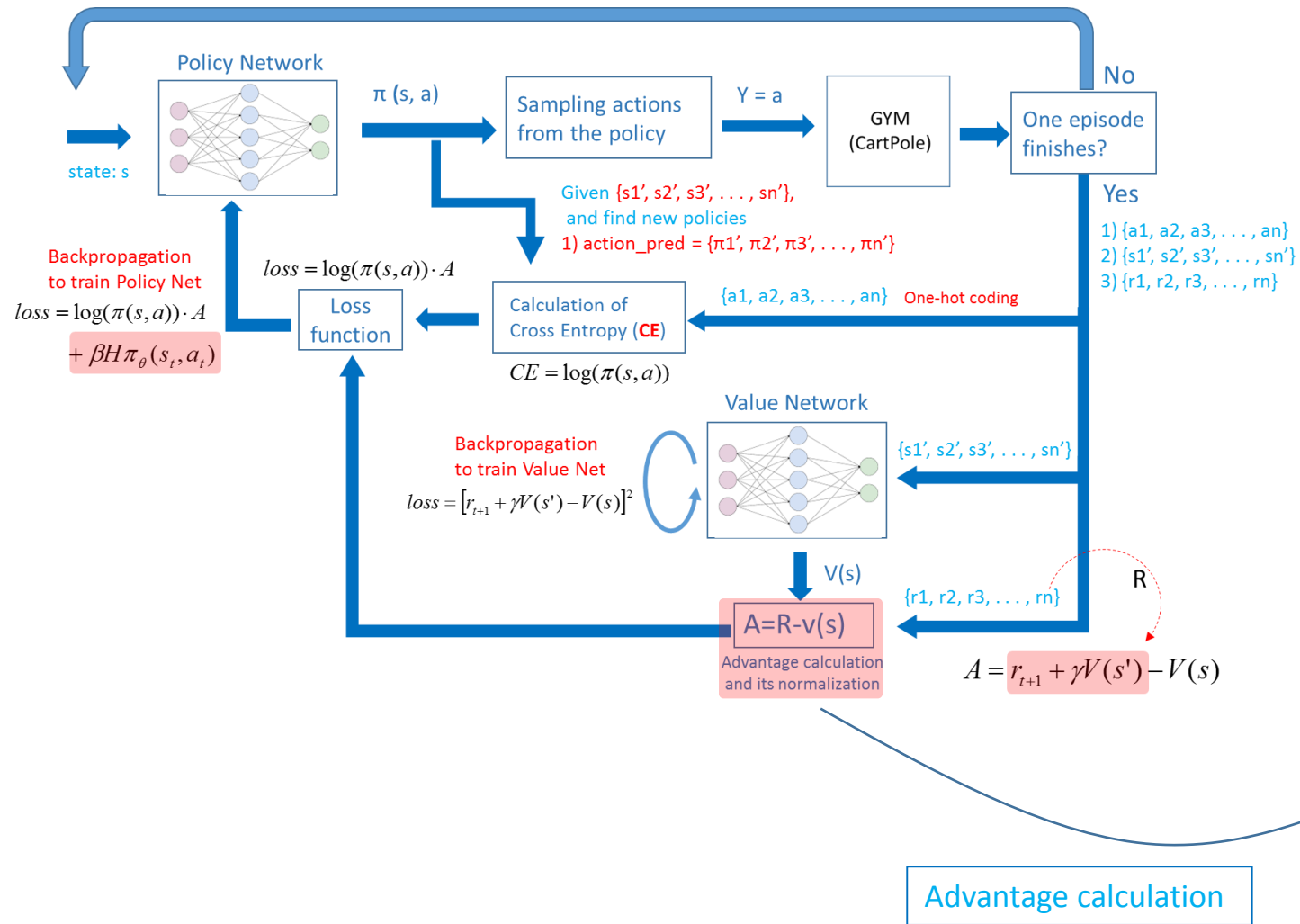
    # 2. Get 'V's
    feed = {
        self.model.S: S
    }
    V = self.sess.run(self.model.V, feed_dict=feed)

    # 3. Get Advantage values, A = R - V
    ADV = R - V
    ADV = (ADV - np.mean(ADV)) / (np.std(ADV) + 1e-8)

    # 4. Perform gradient descents
    feed = {
        self.model.S: S,
        self.model.A: A,
        self.model.ADV: ADV,
        self.model.R: R
    }

    self.sess.run(self.model.train_op, feed_dict=feed)
```





```
def train(self, S, A, R):
    """ Train the actor critic networks

    (1) Compute discounted rewards R
    (2) Compute advantage values A = R - V
    (3) Perform gradients updates

    """

    def discount_rewards(r, gamma=0.99):
        """ take 1D float array of rewards and compute discounted rewards """
        discounted_r = np.zeros_like(r, dtype=np.float32)
        running_add = 0
        for t in reversed(range(len(r))):
            running_add = running_add * gamma + r[t]
            discounted_r[t] = running_add

        return discounted_r

    # 1. Get discounted 'R's
    R = discount_rewards(R)

    # 2. Get 'V's
    feed = {
        self.model.S: S
    }
    V = self.sess.run(self.model.V, feed_dict=feed)

    # 3. Get Advantage values, A = R - V
    ADV = R - V
    ADV = (ADV - np.mean(ADV)) / (np.std(ADV) + 1e-8)

    # 4. Perform gradient descents
    feed = {
        self.model.S: S,
        self.model.A: A,
        self.model.ADV: ADV,
        self.model.R: R
    }

    self.sess.run(self.model.train_op, feed_dict=feed)
```

Advantage calculation



## Calculation of cross entropy

$\pi$  : policy

$[\pi(s_1', a_0), \pi(s_1', a_1)]$   
 $[\pi(s_2', a_0), \pi(s_2', a_1)]$   
 $[\pi(s_3', a_0), \pi(s_3', a_1)]$   
 $\vdots$   
 $\vdots$   
 $\vdots$   
 $[\pi(s_n', a_0), \pi(s_n', a_1)]$

A: action  
one-hot coding

$[[0, 1]]$   
 $[1, 0]$   
 $[1, 0]$   
 $\vdots$   
 $\vdots$   
 $\vdots$   
 $[0, 1]$

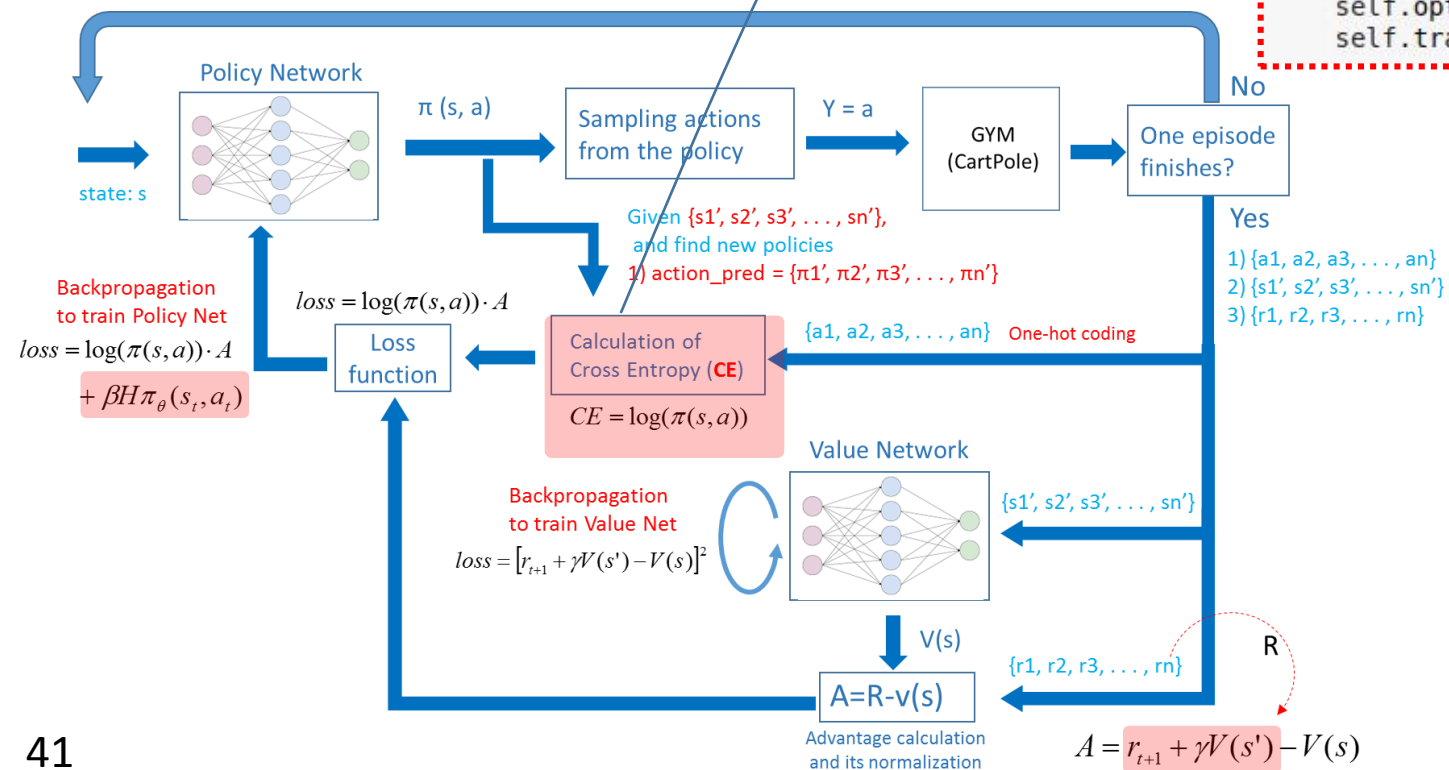
```
def _create_op(self):
    # output shape: [None]
    policy_gain = tf.reduce_sum(self.P * self.A, 1)

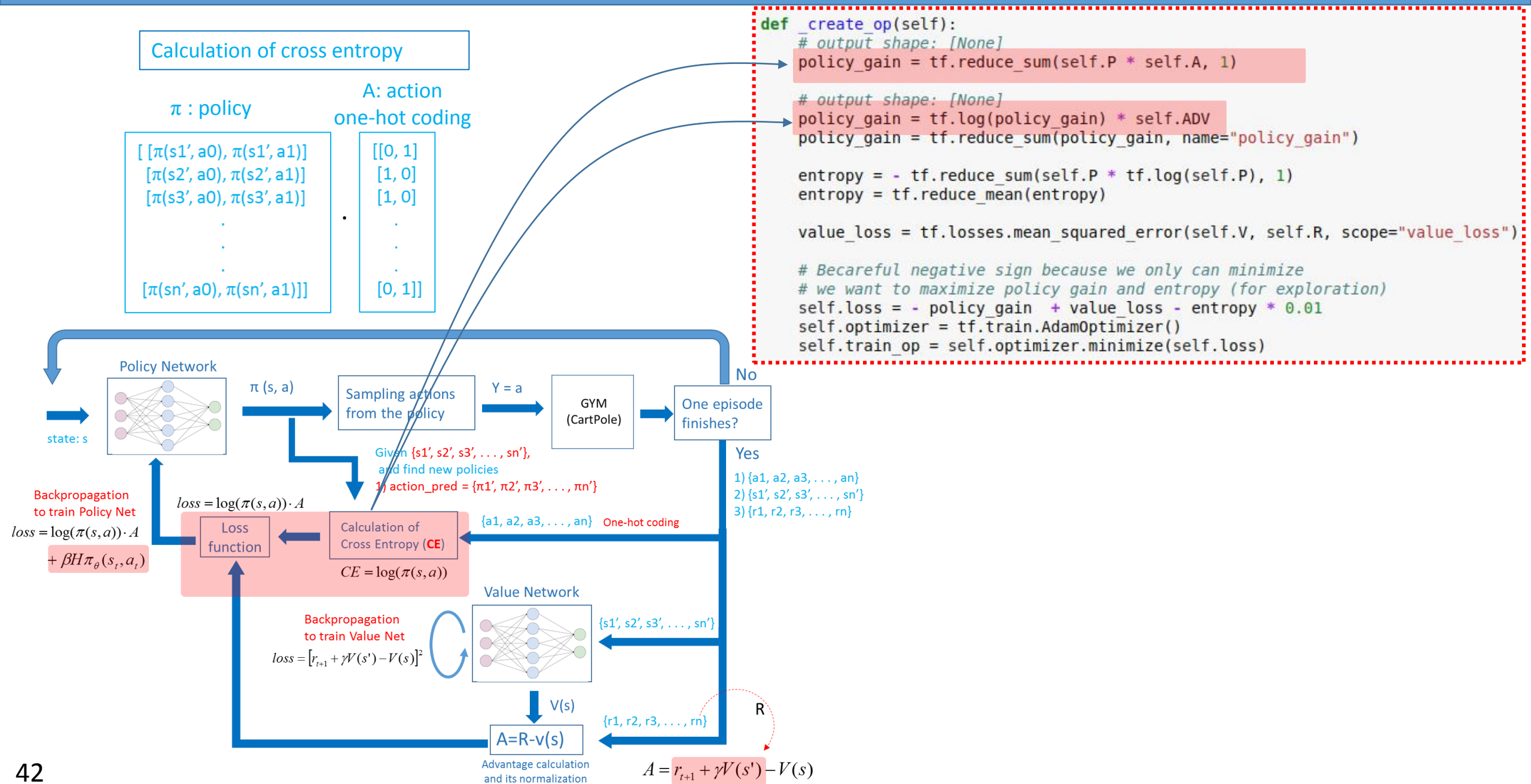
    # output shape: [None]
    policy_gain = tf.log(policy_gain) * self.ADV
    policy_gain = tf.reduce_sum(policy_gain, name="policy_gain")

    entropy = - tf.reduce_sum(self.P * tf.log(self.P), 1)
    entropy = tf.reduce_mean(entropy)

    value_loss = tf.losses.mean_squared_error(self.V, self.R, scope="value_loss")

    # Becareful negative sign because we only can minimize
    # we want to maximize policy gain and entropy (for exploration)
    self.loss = - policy_gain + value_loss - entropy * 0.01
    self.optimizer = tf.train.AdamOptimizer()
    self.train_op = self.optimizer.minimize(self.loss)
```





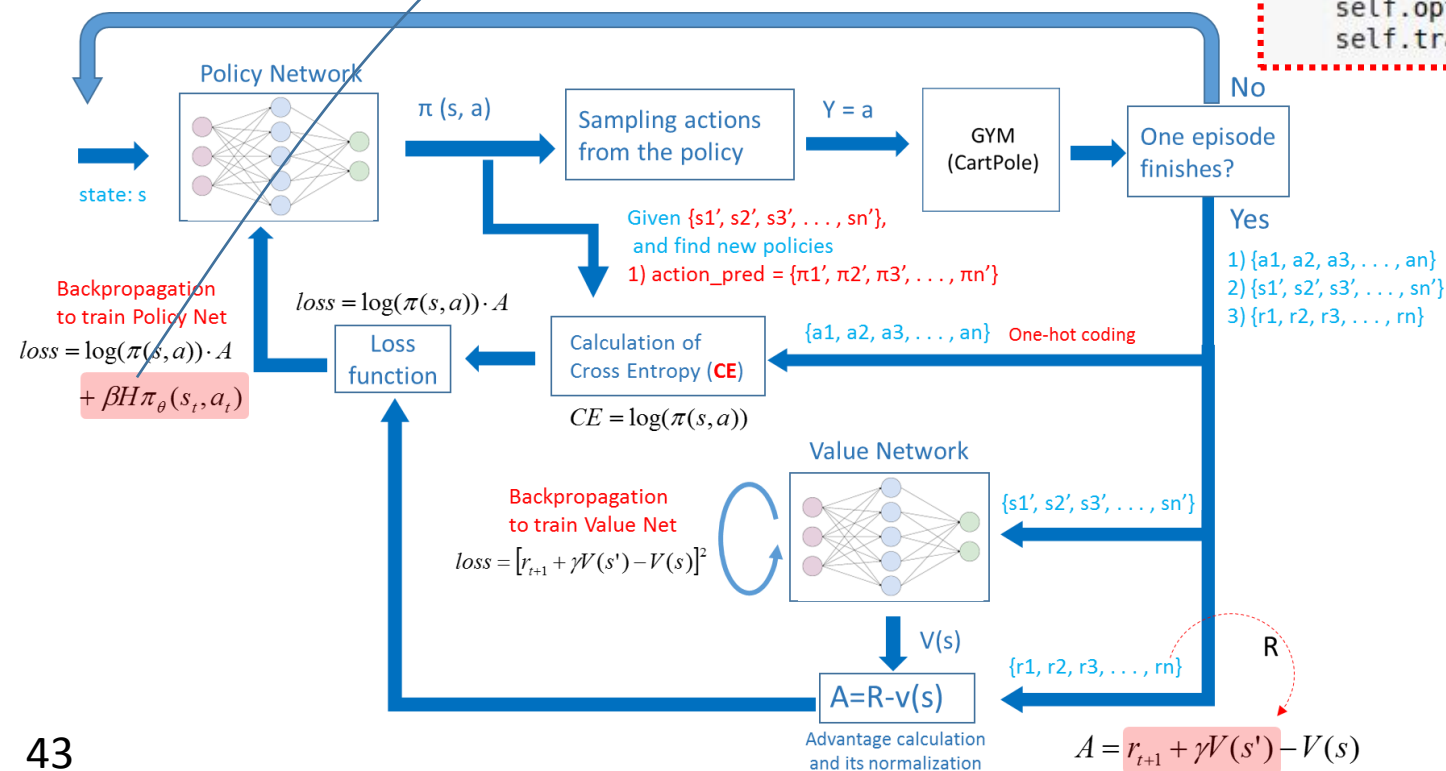
```
def _create_op(self):
    # output shape: [None]
    policy_gain = tf.reduce_sum(self.P * self.A, 1)

    # output shape: [None]
    policy_gain = tf.log(policy_gain) * self.ADV
    policy_gain = tf.reduce_sum(policy_gain, name="policy_gain")

    entropy = - tf.reduce_sum(self.P * tf.log(self.P), 1)
    entropy = tf.reduce_mean(entropy)

    value_loss = tf.losses.mean_squared_error(self.V, self.R, scope="value_loss")

    # Becareful negative sign because we only can minimize
    # we want to maximize policy gain and entropy (for exploration)
    self.loss = - policy_gain + value_loss - entropy * 0.01
    self.optimizer = tf.train.AdamOptimizer()
    self.train_op = self.optimizer.minimize(self.loss)
```



- Entropy regularization term
- This term tries to uniformize the probability distribution of actions defined in the first term.
  - Entropy is maximized when all actions from the policy  $\pi$  are same.
  - It aims to occur all action with equal probability (exploration)

# AC Code: Loss function

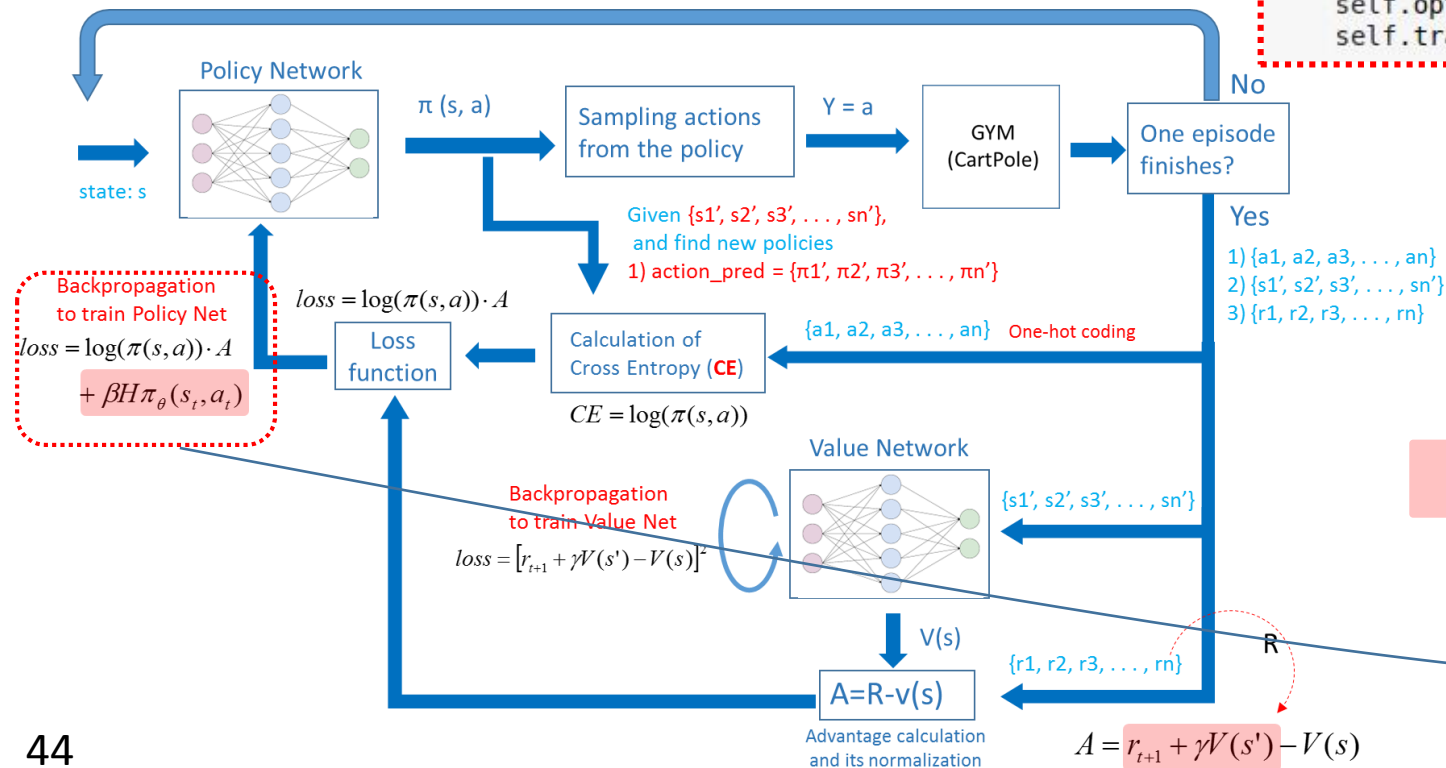
```
def _create_op(self):
    # output shape: [None]
    policy_gain = tf.reduce_sum(self.P * self.A, 1)

    # output shape: [None]
    policy_gain = tf.log(policy_gain) * self.ADV
    policy_gain = tf.reduce_sum(policy_gain, name="policy_gain")

    entropy = - tf.reduce_sum(self.P * tf.log(self.P), 1)
    entropy = tf.reduce_mean(entropy)

    value_loss = tf.losses.mean_squared_error(self.V, self.R, scope="value_loss")

    # Becareful negative sign because we only can minimize
    # we want to maximize policy gain and entropy (for exploration)
    self.loss = - policy_gain + value_loss - entropy * 0.01
    self.optimizer = tf.train.AdamOptimizer()
    self.train_op = self.optimizer.minimize(self.loss)
```



Loss function to train policy Network

$$L(\theta) = \log(\pi_{\theta}(s_t, a_t)) \cdot (r_{t+1} + \gamma V(s_{t+1}) - V(s_t)) + \beta H \pi_{\theta}(s_t, a_t)$$

$$A = r_{t+1} + \gamma V(s') - V(s)$$

Backup Slides



## Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

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  $\epsilon$  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  $\theta$

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

**End For**

**End For**

- Training the main Q-Net and target  $\hat{Q}$ -Net given minibatch input

```
# store the previous observations in replay memory
replay_buffer = deque(maxlen=REPLAY_MEMORY)
```

```
last_100_game_reward = deque(maxlen=100)
```

```
with tf.Session() as sess:
```

```
def replay_train(mainDQN: dqn.DQN, targetDQN: dqn.DQN):
```

```
    """Trains mainDQN with target Q values given by targetDQN"""
```

```
    Args:
```

```
        mainDQN (dqn.DQN): Main DQN that will be trained
```

```
        targetDQN (dqn.DQN): Target DQN that will provide target Q values
```

```
        train_batch (list): Minibatch of replay memory
```

```
            Each element is (s, a, r, s', done)
```

```
            [(state, action, reward, next_state, done)]
```

```
    Returns:
```

```
        float: After updating 'mainDQN', it returns the current value of the last_100_game_reward
```

```
    states = np.vstack([x[0] for x in train_batch])
```

```
    actions = np.array([x[1] for x in train_batch])
```

```
    rewards = np.array([x[2] for x in train_batch])
```

```
    next_states = np.vstack([x[3] for x in train_batch])
```

```
    done = np.array([x[4] for x in train_batch])
```

```
    X = states
```

```
    Q_target = rewards + DISCOUNT_RATE * np.max(np.vstack([targetDQN.predict(next_states), targetDQN.predict(states)]), axis=1)
```

```
    y = mainDQN.predict(states)
```

```
    y[np.arange(len(X)), actions] = Q_target
```

```
    # Train our network using target and predicted values
```

```
    return mainDQN.update(X, y)
```

```
Q = [10, 20, 30]
```

```
y = np.array([[1,2], [4,7], [5,3]])
```

```
print ("Before: ", "\n", y)
```

```
action = [1,0,1]
```

```
y[np.arange(3), action] = Q
```

```
print ("After: ", "\n", y)
```

Before:

```
[[1 2]
```

```
[4 7]
```

```
[5 3]]
```

After:

```
[[ 1 10]
```

```
[20  7]
```

```
[ 5 30]]
```

- 3 batch samples

- Each sample: Q(R), Q(L)

- Insert target Q value into Y at location; a=[1,0,1]

- Q network needs to be trained to produce this value

Tricky part!

```
reward = -1
```

```
# Save the experience to our buffer
```

```
replay_buffer.append((state, action, reward, next_state, done))
```

```
if len(replay_buffer) > BATCH_SIZE:
```

```
    minibatch = random.sample(replay_buffer, BATCH_SIZE)
```

```
    loss, _ = replay_train(mainDQN, targetDQN, minibatch)
```

```
if step_count % TARGET_UPDATE_FREQUENCY == 0:
```

```
    sess.run(copy_ops)
```

```
state = next_state
```

```
step_count += 1
```

```
print("Episode: {} steps: {}".format(episode, step_count))
```