# 7.강 Deep Q-learning
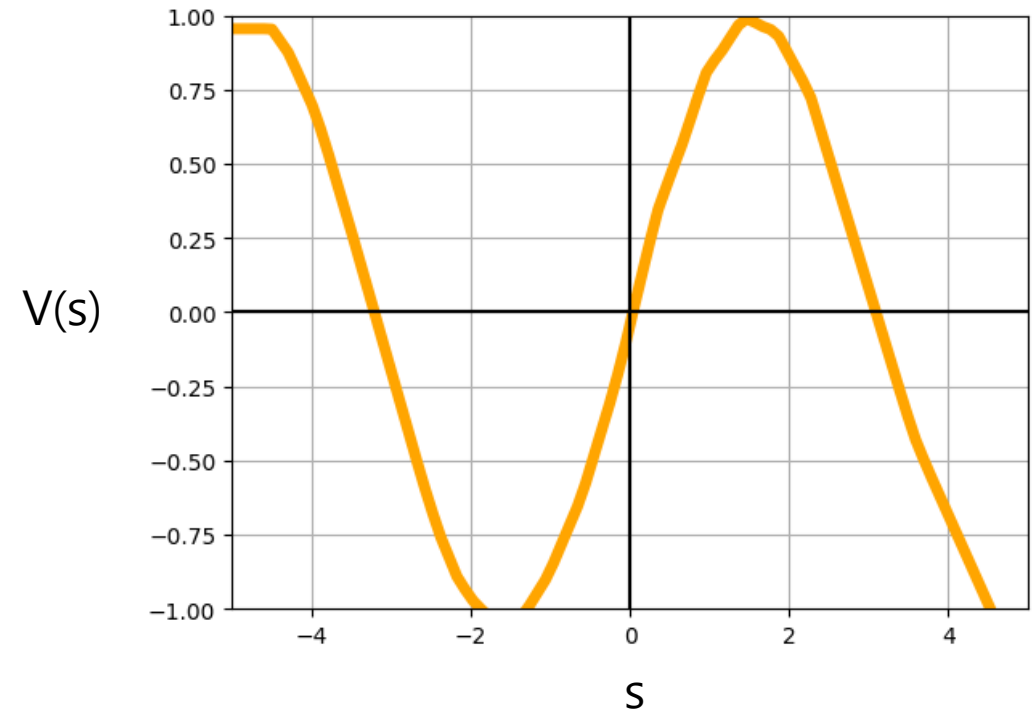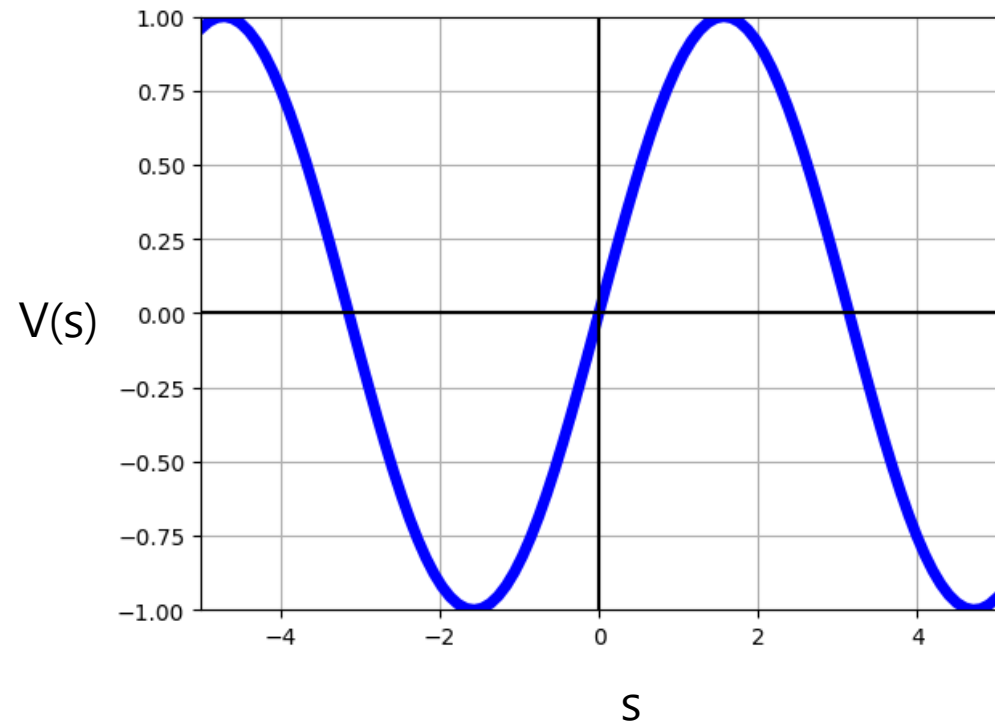
# Contents

- Neural Network
- Pytorch
- Deep SARSA
- Code Exercise
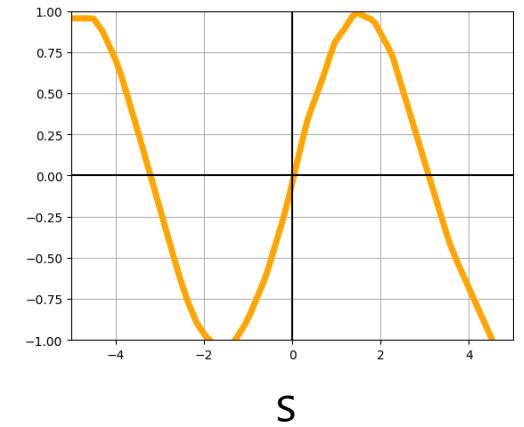
# Neural Network

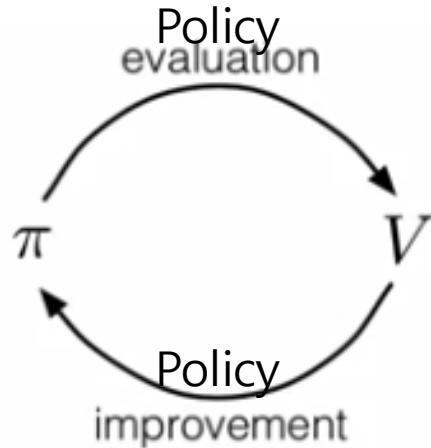# Function approximators

V(s)

V(s)

# Function approximators

- How do we observe the value function?
    - The agent learns based on experience.
    - The functions v*(s) and q*(s ,a) are not know in advance



$$f_1(s|w) \rightarrow f_2(s|w) \rightarrow \cdots . \rightarrow f_n(s|w) \approx v*(s)$$

# Neural Networks

- Computing system inspired by the biological neural networks that constitute our brain

- They server multiple purposes, including function approximation: $\hat{y} = f(x|w)$

- Mathematical function typically consisting of a weighted sum of inputs and a activation/transfer function $Output = \varphi(\sum_{i=1}^{n} w_i x_i + b)$



(a) Biological neuron

(b) Artificial neuron

# Neural networks

- Input vector x = [x1, x2, x3]

- Connection matrix:

$$W1 = \begin{bmatrix} w11 & w12 & w13 & w14 \\ w21 & w22 & w23 & w24 \\ w31 & w32 & w33 & w34 \end{bmatrix}$$

By changing its parameters W1, we can modify it to approximate the function we are interested in

- Output vector:
$$H = [\varphi(\sum_{i=1}^{n} w_i x_i + b), ..., \varphi(\sum_{i=1}^{n} w_i x_i + b)]$$

# Neural Networks

- Networks that do not have cycles are known as feedforward NN. Signals always propagate forward

- The neuron receives inputs, process & aggregate those inputs, and either inhibits or amplifies before passing the signal to the next layer

input layer

hidden layer 1    hidden layer 2

output layer

$$h_1 = \sigma(W_1 x) \qquad h_2 = \sigma(W_2 h_1) \qquad y = W_3 h_2$$

# Neural networks

- Activation functions

Relu()

Sigmoid()

Tanh()

$\varphi(Vk)$

$V_k$

$V_k$

$V_k$

# Gradient Descent

- Given some loss function: $L(\vec{x}, \vec{y}) = \|2\vec{x} + 2\vec{y}\|$
- Update rules for the parameters: $w_{t+1} = w_t - \alpha \nabla \hat{L}(w)$
- Gradient vector: $\nabla \hat{L}(w) = [\frac{\partial L}{\partial w1}, \frac{\partial L}{\partial w2}, \ldots, \frac{\partial L}{\partial wn}]$
- Computed using the backpropagation algorithm
- $\nabla \hat{L}(w)$ points to the direction of maximum growth of $\nabla \hat{L}(w)$
- $\alpha$ is the size of the step we take in the opposite direction to $\nabla \hat{L}(w)$

# Backpropagation

- Computed using the backpropagation algorithm
- We want to evaluate partial derivative: $\frac{\partial L}{\partial \vec{x}}$ and $\frac{\partial L}{\partial \vec{y}}$

```
shape = (3, )
x = torch.tensor([1., 2, 3], requires_grad=True)
y = torch.ones(shape, requires_grad=True)

loss = ((2 * x + y)**2).sum()
loss.backward()
print(x.grad)
print(y.grad)
```

```
tensor([24., 40., 56.])
tensor([12., 20., 28.])
```



```
loss = ((2 * x + y)**2).sum()
print(loss)
```

```
tensor(83., grad_fn=<SumBackward0>)
```

# Cost function

- Mean squared error:

$$L(w) = \frac{1}{N} \sum_{i=0}^{N} [y - \hat{y}]^2$$

- For our neural network to estimate q(s, a) as well as possible, we will minimize the observed squared errors

$$\hat{L}(w) = \frac{1}{N} \sum_{i=0}^{N} [R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1} | w) - \hat{q}(S_t, A_t | w)]^2$$

- Target value:                                        Estimated value:

$$R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1} | w) \qquad\qquad \hat{y} = \hat{q}(S_t, A_t | w)$$

# Neural network optimization

- For our neural network to estimate q(s, a) as well as possible, we will minimize the observed squared errors

$$\hat{L}(\theta) = \frac{1}{N}\sum_{i=0}^{N}[R_i + \gamma \hat{q}(S_i', A_i'|\theta_{targ}) - \hat{q}(S_i, A_i|\theta)]^2$$

- Target value: a value towards which we want to push the estimates

$$R_i + \gamma \hat{q}(S_i', A_i'|\theta_{targ})$$

- Estimate of the q-value of a state-action pair

$$\hat{q}(S_i, A_i|\theta)$$

# Pytorch

# Numpy & PyTorch

**NumPy**

- Fast CPU implementations
- CPU-only
- No autodiff
- Imperative

**PyTorch**

- Fast CPU implementations
- Allows GPU
- Supports autodiff
- Imperative

**Other features include:**

- Datasets and dataloading
- Common neural network operations
- Built-in optimizers (Adam, SGD, …)

# Multidimensional Indexing



A.shape == (3, 5, 4)

# Shape Operations

NumPy

```python
A = np.random.normal(size=(10, 15))

# Indexing with newaxis/None
# adds an axis with size 1
A[np.newaxis] # -> shape (1, 10, 15)

# Squeeze removes a axis with size 1
A[np.newaxis].squeeze(0) # -> shape (10, 15)

# Transpose switches out axes.
A.transpose((1, 0)) # -> shape (15, 10)

# !!! BE CAREFUL WITH RESHAPE !!!
A.reshape(15, 10)    # -> shape (15, 10)
A.reshape(3, 25, -1) # -> shape (3, 25, 2)
```

PyTorch

```python
A = torch.randn((10, 15))

# Indexing with None
# adds an axis with size 1
A[None] # -> shape (1, 10, 15)

# Squeeze removes a axis with size 1
A[None].squeeze(0) # -> shape (10, 15)

# Permute switches out axes.
A.permute((1, 0)) # -> shape (15, 10)

# !!! BE CAREFUL WITH VIEW !!!
A.view(15, 10)    # -> shape (15, 10)
A.view(3, 25, -1) # -> shape (3, 25, 2)
```

# Device Management

- Numpy: all arrays live on the CPU's RAM
- Torch: tensors can either live on CPU or GPU memory
  - Move to GPU with .to("cuda")/.cuda()
  - Move to CPU with .to("cpu")/.cpu()

## YOU CANNOT PERFORM OPERATIONS BETWEEN TENSORS ON DIFFERENT DEVICES!

# Computing Gradients

```python
P = torch.randn((1024, 1024))
print(P.requires_grad) # -> False
P = torch.randn((1024, 1024), requires_grad=True)
b = torch.randn((1024,), requires_grad=True)
print(P.grad) # -> None

x = torch.randn((32, 1024))
y = torch.nn.relu(x @ P + b)

target = 3
loss = torch.mean((y - target) ** 2).detach()
```

# Training Loop

**REMEMBER THIS!**

```python
net = (...).to("cuda")
dataset = ...
dataloader = ..
optimizer = ...
loss_ fn = ..
for epoch in range(num_epochs):
  # Training..
  net.train()
  for data, target in dataloader:
    data = torch.from_numpy(data).float().cuda()
    target = torch.from_numpy(data).float().cuda()

    prediction = net(data)
    loss = loss_fn(prediction, target)

    optimizer.zero_grad()
    loss.backward()
    optimizer.step()

  net.eval()
  # Do evaluation..
```

# Converting Numpy / PyTorch

**Numpy -> PyTorch:**

```
torch.from_numpy(numpy_array).float()
```

**PyTorch -> Numpy:**

- (If `requires_grad`) Get a copy without graph with `.detach()`
- (If on GPU) Move to CPU with `.to("cpu")`/`.cpu()`
- Convert to numpy with `.numpy`

**All together:**

```
torch_tensor.detach().cpu().numpy()
```

# Custom networks

```python
import torch.nn as nn

class SingleLayerNetwork(nn.Module):
    def __init__(self, in_dim: int, out_dim: int, hidden_dim: int):
        super().__init__() # <- Don't forget this!
        self.net = nn.Sequential(
            nn.Module(in_dim, hidden_dim),
            nn.ReLU(),
            nn.Module(hidden_dim, out_dim),
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        return self.net(x)

batch_size = 256
my_net = SingleLayerNetwork(2, 32, 1).to("cuda")
output = my_net(torch.randn(size=(batch_size, 2)).cuda())
```

- nn.Module represents the building blocks of a computation graph.
  - For example, in typical pytorch code, each convolution block is its own module, each fully connected block is a module, and the whole network itself is also a module.

- Modules can contain modules within them. All the classes inside of `torch.nn` are instances `nn.Modules`.

# Custom networks

```python
import torch.nn as nn

class SingleLayerNetwork(nn.Module):
    def __init__(self, in_dim: int, out_dim: int, hidden_dim: int):
        super().__init__() # <- Don't forget this!
        self.net = nn.Sequential(
            nn.Module(in_dim, hidden_dim),
            nn.ReLU(),
            nn.Module(hidden_dim, out_dim),
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        return self.net(x)

batch_size = 256
my_net = SingleLayerNetwork(2, 32, 1).to("cuda")
output = my_net(torch.randn(size=(batch_size, 2)).cuda())
```

- Prefer `net()` over `net.forward()`
- Everything (network and its inputs) on the same device!!!

# Torch Best Practices

- When in doubt, **assert** is your friend

```
assert x.shape == (B, N), \
        f"Expected shape ({B, N}) but got {x.shape}"
```

- Be extra careful with `.reshape/.view`
  - If you use it, assert before and after
  - Only use it to collapse/expand a single dim
  - In Torch, prefer `.flatten()/.permute()/.unflatten()`

- If you do some complicated operation, test it!
  - Compare to a pure Python implementation

# Torch Best Practices

- Don't mix numpy and Torch code
  - Understand the boundaries between the two
  - Make sure to cast 64-bit numpy arrays to 32 bits
  - torch.Tensor only in nn.Module!
- Training loop will always look the same
  - Load batch, compute loss
  - .zero_grad(), .backward(), .step()

# Neural network optimization

- Mean squared error:

$$L(\theta) = \frac{1}{N} \sum_{i=1}^{N} [y_i - \hat{y}_i]^2$$

- We want to minimize the square of the errors of the neural network estimates

# Neural network optimization

- We calculate the gradient vector of the cost function with respect to the $\theta$ parameters:

$$\nabla L(\theta) = \left[\frac{\partial L}{\partial \theta 1}, \frac{\partial L}{\partial \theta 2}, \ldots, \frac{\partial L}{\partial \theta n}\right]$$

- With the gradient vector, we will make a SGD step:

$$\theta \leftarrow \theta - \alpha \nabla \hat{L}(\theta)$$

# Neural Net. Architecture for V, Q

- St vector input → NN → V scalar output

- St, At vector input → NN → Q scalar output
  - Continuous case

- St vector input → Q vector output
  - Discrete action space only
  - Output size is $|A|$

# Neural Net. Architecture for policy, $\pi$

- St input $\rightarrow$ NN $\rightarrow$ vector output
  - Discrete action space case
  - Output size is $|A|$
  - SOFTMAX turns the output into prob. (sum of prob. is 1)


- St input $\rightarrow$ NN $\rightarrow$ $\mu_\theta(S_t), \delta_\theta(St)$ output
  - Continuous action space case
  - Represented with Gaussian Distribution

# Deep SARSA

# Neural network optimization

$$L(\theta) = \frac{1}{|K|}\sum_{i=0}^{|K|}[R_i + \gamma\hat{q}(S_i', A_i'|\theta_{targ}) - \hat{q}(S_i, A_i|\theta)]^2$$

- Target is the value towards which we want to push the estimates.

$$R_i + \gamma\hat{q}(S_i', A_i'|\theta_{targ})$$

- Estimate is the estimate of the q-value of a state-action pair
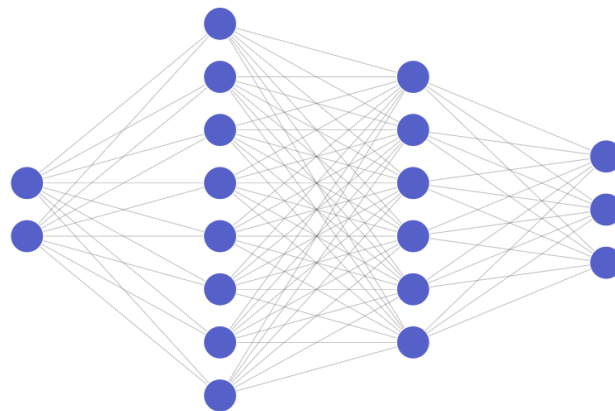
$$\hat{q}(S_i, A_i|\theta)$$

# Target network

Bootstrapping  **+**  Function approximator



$$y_i = R_i + \gamma \hat{q}(S_i', A_i'|\theta_{targ})$$

# Target network

- When a value is changed, nearby values will also be affected.

- By modifying a $\hat{q}(S_i, Ai|\theta)$ estimate we also modify its $\hat{q}(S_i', A_i'|\theta_{targ})$ target

- For the learning process to be stable, the target must also be stable

- power of neural networks



q(s, a)

- - - - Before update
———— After update

s

# Target network

- We make a copy of the neural network to calculate the targets.

$$\theta_{targ} \leftarrow \theta$$

- This neural network does not change with SGD. Its $\theta$ parameters remain the same

- The estimated value of $S_i'$, $A_i'$ is calculated with the target network:

$$L(\theta) = \frac{1}{N}\sum_{i=0}^{N}[R_i + \boxed{\gamma \hat{q}(S_i', A_i'|\theta_{targ})} - \hat{q}(S_i, A_i|\theta)]^2$$
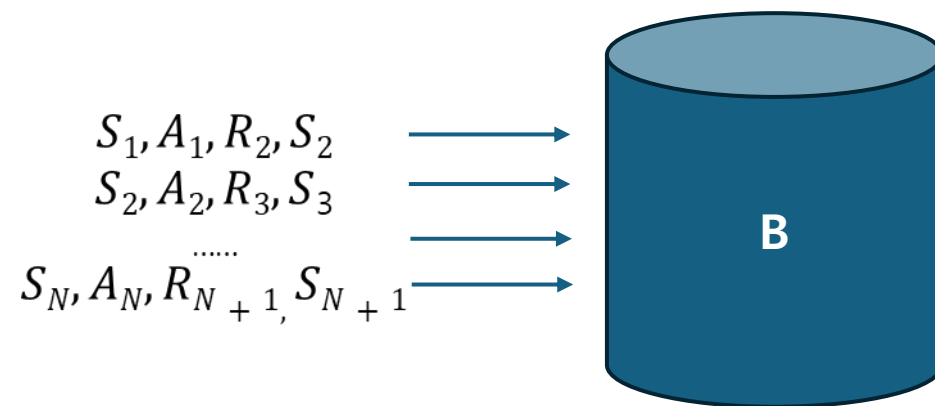
# Neural network optimization

**Algorithm 1** Deep SARSA

1: **Input:** $\alpha$ learning rate, $\epsilon$ random action probability,
2:     $\gamma$ discount factor,
3: Initialize q-value parameters $\theta$ and target parameters $\theta_{targ} \leftarrow \theta$
4: $\pi \leftarrow \epsilon$-greedy policy w.r.t $\hat{q}(s, a|\theta)$
5: Initialize replay buffer $B$
6: **for** episode $\in 1..N$ **do**
7:     Restart environment and observe the initial state $S_0$
8:     **for** $t \in 0..T-1$ **do**
9:         Select action $A_t \sim \pi(S_t)$
10:         Execute action $A_t$ and observe $S_{t+1}, R_{t+1}$
11:         Insert transition $(S_t, A_t, R_{t+1}, S_{t+1})$ into the buffer $B$
12:         $K = (S, A, R, S') \sim B$
13:         Select actions $A' \sim \pi(S')$
14:         Compute loss function over the batch of experiences:

$$L(\theta) = \frac{1}{|K|} \sum_{i=1}^{|K|} [R_i + \gamma \hat{q}(S'_i, A'_i|\theta_{targ}) - \hat{q}(S_i, A_i|\theta)]^2 \qquad (1)$$

15:     **end for**
16:     Every $k$ episodes synchronize $\theta_{targ} \leftarrow \theta$
17: **end for**
18: **Output:** Near optimal policy $\pi$ and q-value approximations $\hat{q}(s, a|\theta)$
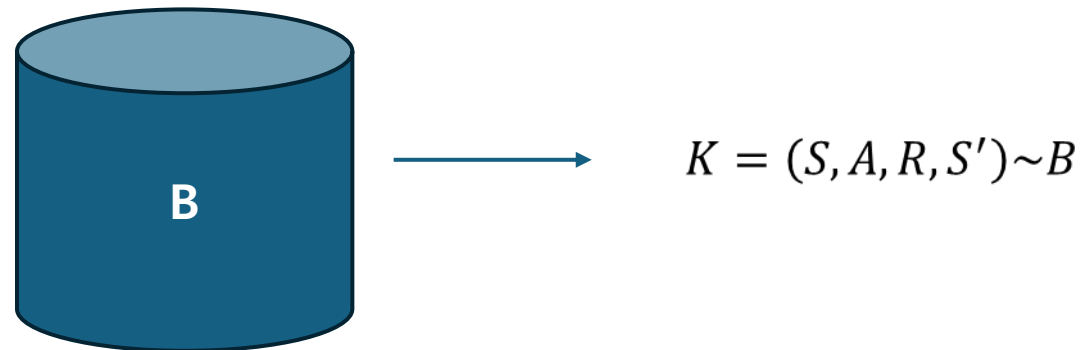
# Experience Replay

- Memory that stores the state transition that the agent experiences
- The memory has a limited size and when it fills up, it replaces old transitions with new ones

$$S_1, A_1, R_2, S_2$$
$$S_2, A_2, R_3, S_3$$
$$......$$
$$S_N, A_N, R_{N+1}, S_{N+1}$$

B

# Experience Replay

- To update the neural network, we randomly chose a batch of transitions from the memory



$$K = (S, A, R, S') \sim B$$

- The batch of transitions obtained from the memory is used to calculate the cost function and update the $\theta$ parameters

- $L(\theta) = \frac{1}{|K|} \sum_{i=0}^{|K|} [R_i + \gamma \hat{q}(S_i', A_i' | \theta_{targ}) - \hat{q}(S_i, A_i | \theta)]^2$

# Neural network optimization

**Algorithm 1** Deep SARSA

1: **Input:** $\alpha$ learning rate, $\epsilon$ random action probability,
2:     $\gamma$ discount factor,
3: Initialize q-value parameters $\theta$ and target parameters $\theta_{targ} \leftarrow \theta$
4: $\pi \leftarrow \epsilon$-greedy policy w.r.t $\hat{q}(s, a | \theta)$
5: Initialize replay buffer $B$
6: **for** episode $\in 1..N$ **do**
7:     Restart environment and observe the initial state $S_0$
8:     **for** $t \in 0..T-1$ **do**
9:         Select action $A_t \sim \pi(S_t)$
10:         Execute action $A_t$ and observe $S_{t+1}, R_{t+1}$
11:         Insert transition $(S_t, A_t, R_{t+1}, S_{t+1})$ into the buffer $B$
12:         $K = (S, A, R, S') \sim B$
13:         Select actions $A' \sim \pi(S')$
14:         Compute loss function over the batch of experiences:

$$L(\theta) = \frac{1}{|K|} \sum_{i=1}^{|K|} [R_i + \gamma \hat{q}(S'_i, A'_i | \theta_{targ}) - \hat{q}(S_i, A_i | \theta)]^2 \tag{1}$$

15:     **end for**
16:     Every $k$ episodes synchronize $\theta_{targ} \leftarrow \theta$
17: **end for**
18: **Output:** Near optimal policy $\pi$ and q-value approximations $\hat{q}(s, a | \theta)$

# Code Ex.

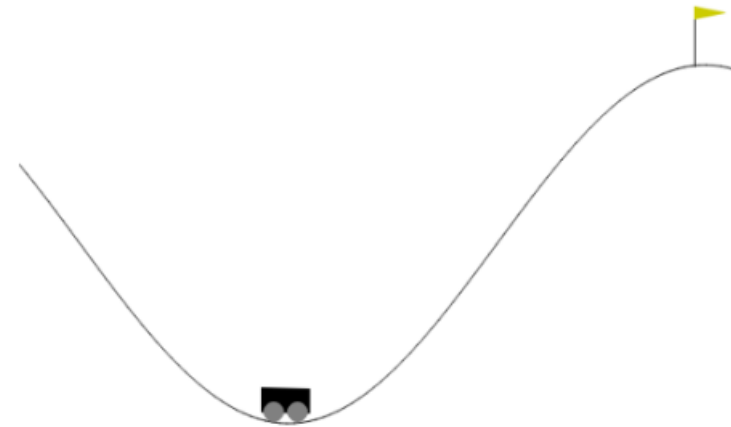• MountainCar: Reach the goal from the bottom of the valley

The state

The observation space consists of the car position $\in [-1.2, 0.6]$ and car velocity $\in [-0.07, 0.07]$

The actions available

The actions available three:

```
0    Accelerate to the left.
1    Don't accelerate.
2    Accelerate to the right.
```

# Code Ex.

- deep_sarsa_colab.ipynb