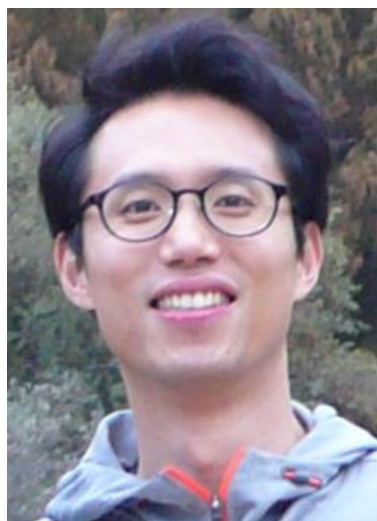


REINFORCE ALGORITHM



Prof. Jae Young Choi

Pattern Recognition and Machine Intelligence Lab. (PMI)

Hankuk University of Foreign Studies

REINFORCE-THEORY

❖ Key Idea

- During learning, actions that resulted in good outcomes should become more probable
- These actions are positively reinforced
- Action probabilities are changed by following the policy gradient
- REINFORCE is known as policy gradient algorithm
- **Policy gradient** is used to modify the policy parameters to maximize the objective

REINFORCE-THEORY

❖ Objective

- Return

$$R_t(\tau) = \sum_{t'=t}^T \gamma^{t'-t} r'_{t'}$$

$$\tau = s_0, a_0, r_0, \dots, s_T, a_T, r_T$$

The return of a trajectory $R_t(\tau)$ is discounted sum of rewards from time step t to the end of a trajectory

REINFORCE-THEORY

❖ Objective

$$J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)] = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \gamma^t r_t \right]$$

The objective $J(\pi_{\theta})$ is the expected return over all complete trajectories generated by an agent

Expectation is calculated over many trajectories sampled from a policy, $\tau \sim \pi_{\theta}$

REINFORCE-THEORY

❖ The Policy π_θ and Objective $J(\pi_\theta)$

- The policy provides a way for an agent to act
- The objective provides a target to maximize

❖ Solving the following problem

$$\max_{\theta} J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)]$$

❖ To maximize the objective, we perform gradient ascent

$$\theta \leftarrow \theta + \alpha \underline{\nabla_{\theta} J(\pi_\theta)}$$

Compute gradient and use it to update the parameters

REINFORCE-THEORY

❖ The policy gradient $\nabla_{\theta} J(\pi_{\theta})$

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \overbrace{R_t(\tau)}^{\text{The return for trajectory } \tau \sim \pi_{\theta}} \underbrace{\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)}_{\text{The probability of the action taken by the agent at time step 't'}} \right]$$

- The action is sampled from the policy $a_t \sim \pi_{\theta}(s_t)$
- The statement of policy gradient
→ Expected sum of the gradients of the log probabilities of the action a_t multiplied by the corresponding return $R_t(\tau)$

REINFORCE-THEORY

❖ Interpreting policy gradient

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T \boxed{R_t(\tau)} \nabla_{\theta} \log \pi_{\theta}(\underline{a_t | s_t}) \right]$$

The return for trajectory $\tau \sim \pi_{\theta}$

The probability of the action taken by the agent at time step 't'

- If the return $R_t(\tau) > 0$, then the probability of the action $\pi_{\theta}(a_t | s_t)$ is increased.
- If the return $R_t(\tau) < 0$, then the probability of the action $\pi_{\theta}(a_t | s_t)$ is decreased.
- Over the course of many updates, the policy will learn to produce actions which result in high $R_t(\tau)$

REINFORCE-THEORY

❖ Policy Gradient Derivation

$$\nabla_{\theta} J(\pi_{\theta}) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)]$$

- We cannot differentiate $R(\tau) = \sum_{t=0}^T \gamma^t r_t$ with respect to θ
 - Rewards r_t are generated by an unknown reward function $\mathcal{R}(s_t, a_t, s_{t+1})$
- The only way is by changing the state and action distributions which, in turn, change the rewards received by an agent
 - Transform $\nabla_{\theta} J(\pi_{\theta}) = \nabla_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau)]$ into a form where we can take a gradient with respect to θ

REINFORCE-THEORY

❖ Policy Gradient Derivation

$$\begin{aligned} & \nabla_{\theta} \overbrace{\mathbb{E}_{x \sim p(x|\theta)} [f(x)]}^{\text{Expectation of function } f(x)} \\ & \quad \underbrace{\mathbb{E}_{x \sim p(x|\theta)}}_{\text{Parameterized probability distribution}} \underbrace{[f(x)]}_{\text{Function}} \\ &= \nabla_{\theta} \int dx f(x) p(x|\theta) && (\text{definition of expectation}) \\ &= \int dx \nabla_{\theta} (p(x|\theta) f(x)) && (\text{bring in } \nabla_{\theta}) \\ &= \int dx (f(x) \nabla_{\theta} p(x|\theta) + p(x|\theta) \nabla_{\theta} f(x)) && (\text{chain-rule}) \\ &= \int dx f(x) \nabla_{\theta} p(x|\theta) && (\nabla_{\theta} f(x) = 0) \\ &= \int dx f(x) p(x|\theta) \frac{\nabla_{\theta} p(x|\theta)}{p(x|\theta)} && \left(\text{multiply } \frac{p(x|\theta)}{p(x|\theta)} \right) \\ &= \int dx f(x) p(x|\theta) \nabla_{\theta} \log p(x|\theta) && \left\{ \begin{array}{l} \boxed{\nabla_{\theta} \log p(x|\theta) = \frac{\nabla_{\theta} p(x|\theta)}{p(x|\theta)}} \\ (\text{definition of expectation}) \end{array} \right. \\ &= \mathbb{E}_x [f(x) \nabla_{\theta} \log p(x|\theta)] \end{aligned}$$

- $f(x)$ is a black-box function which cannot be integrated
- To deal with it, we convert the equation into an expectation
→ It can be estimated through sampling

REINFORCE-THEORY

❖ Policy Gradient Derivation

$$\nabla_{\theta} \mathbb{E}_{x \sim p(x|\theta)} [f(x)] = \mathbb{E}_x [f(x) \nabla_{\theta} \log p(x | \theta)]$$

Gradient of an expectation
= Expectation of the gradient of log probability
multiplied by the original function

REINFORCE-THEORY

❖ Policy Gradient Derivation

By replacing $x = \tau, f(x) = R(\tau), p(x | \theta) = p(\tau | \theta)$

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) \nabla_{\theta} \log p(\tau | \theta)]$$

REINFORCE-THEORY

❖ Policy Gradient Derivation

By replacing $x = \tau, f(x) = R(\tau), p(x | \theta) = p(\tau | \theta)$

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) \nabla_{\theta} \log p(\tau | \theta)]$$

The term $p(\tau | \theta)$ should be related to the policy $\pi_{\theta}(a_t | s_t)$ (we can control over)

It needs to be expanded further in the next slides

REINFORCE-THEORY

❖ Policy Gradient Derivation

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) \nabla_{\theta} \log p(\tau \mid \theta)]$$

$$p(\tau \mid \theta) = \prod_{t \geq 0} p(s_{t+1} \mid s_t, a_t) \pi_{\theta}(a_t \mid s_t)$$

$$\log p(\tau \mid \theta) = \log \prod_{t \geq 0} p(s_{t+1} \mid s_t, a_t) \pi_{\theta}(a_t \mid s_t)$$

$$\log p(\tau \mid \theta) = \sum_{t \geq 0} (\log p(s_{t+1} \mid s_t, a_t) + \log \pi_{\theta}(a_t \mid s_t))$$

$$\nabla_{\theta} \log p(\tau \mid \theta) = \nabla_{\theta} \sum_{t \geq 0} (\log p(s_{t+1} \mid s_t, a_t) + \log \pi_{\theta}(a_t \mid s_t))$$

$\log p(s_{t+1} \mid s_t, a_t)$ is independent of θ , its gradient is zero

$$\nabla_{\theta} \log p(\tau \mid \theta) = \nabla_{\theta} \sum_{t \geq 0} \log \pi_{\theta}(a_t \mid s_t)$$

REINFORCE-THEORY

❖ Policy Gradient Derivation

- Final formation

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

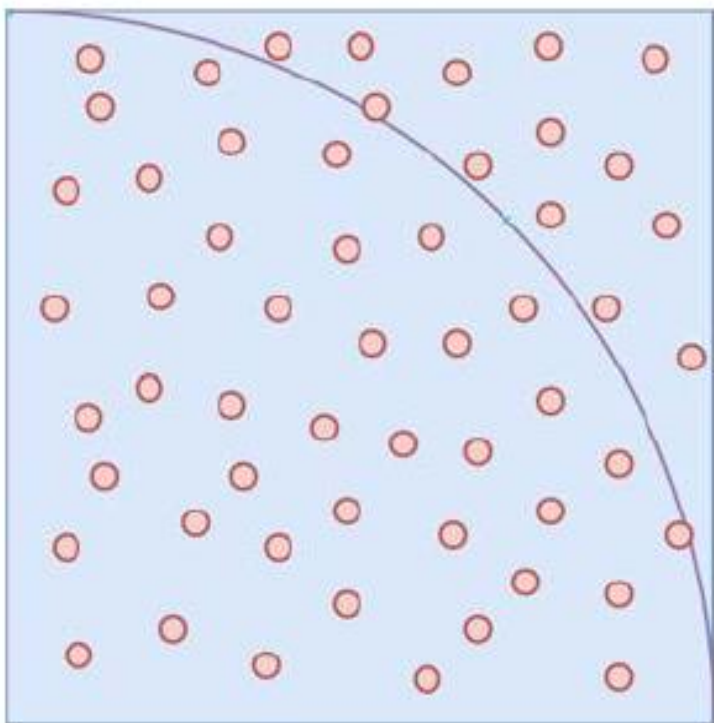
Expressing the probability $p(\tau | \theta)$ in terms of $\pi_{\theta}(a_t | s_t)$

$\nabla_{\theta} J(\pi_{\theta})$ can be estimated easily using a policy network π_{θ} with gradient computation

REINFORCE-THEORY

❖ Implementation of policy gradient

- The REINFORCE algorithm numerically estimates the policy gradient using Monte Carlo sampling
- Use random sampling to generate data used to approximate a function
→ "approximation with random sampling"



$$\frac{\text{area of circle}}{\text{area of square}} = \frac{\pi r^2}{(2r)^2} = \frac{\pi}{4}$$

$$\pi = \frac{\text{No. of samples in a circle}}{\text{No. of total samples}}$$

By iteratively sampling more points and updating the ratio, our estimation will get closer to the precise value

Monte Carlo sampling used to estimate π

REINFORCE-THEORY

❖ Implementation of policy gradient

$$\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^T R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

Expectation means that more trajectories τ are sampled using a policy and averaged, it approaches the actual policy gradient $\nabla_{\theta} J(\pi_{\theta})$

- Monte Carlo estimate over sampled trajectories

$$\nabla_{\theta} J(\pi_{\theta}) \approx \sum_{t=0}^T R_t(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

UNDERSTANDING REINFORCE

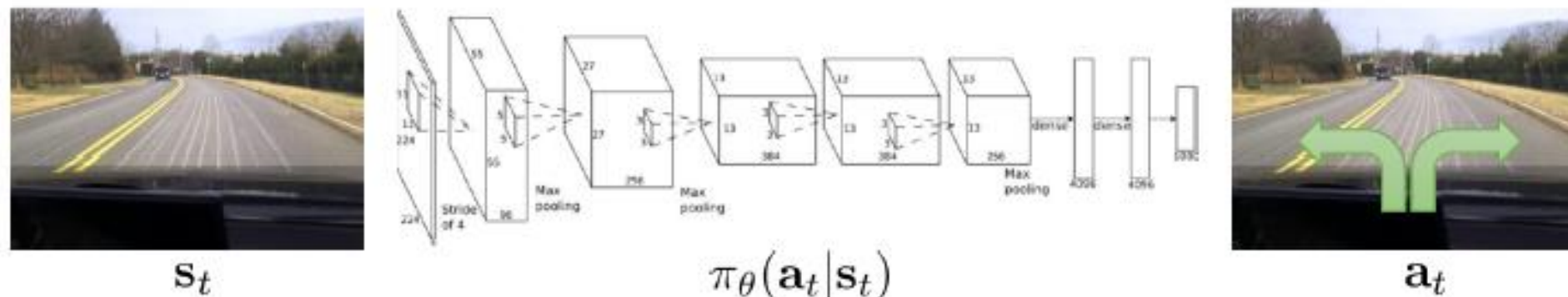
❖ Evaluating the policy gradient

$$\text{recall: } J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\sum_t r(\mathbf{s}_t, \mathbf{a}_t) \right] \approx \frac{1}{N} \sum_i \sum_t r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

$$\nabla_{\theta} J(\theta) = E_{\tau \sim p_{\theta}(\tau)} \left[\left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \right) \left(\sum_{t=1}^T r(\mathbf{s}_t, \mathbf{a}_t) \right) \right]$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

what is this?



UNDERSTANDING REINFORCE

❖ What did we just do?

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^T r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau_i)}_T r(\tau_i)$$
$$\sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$$

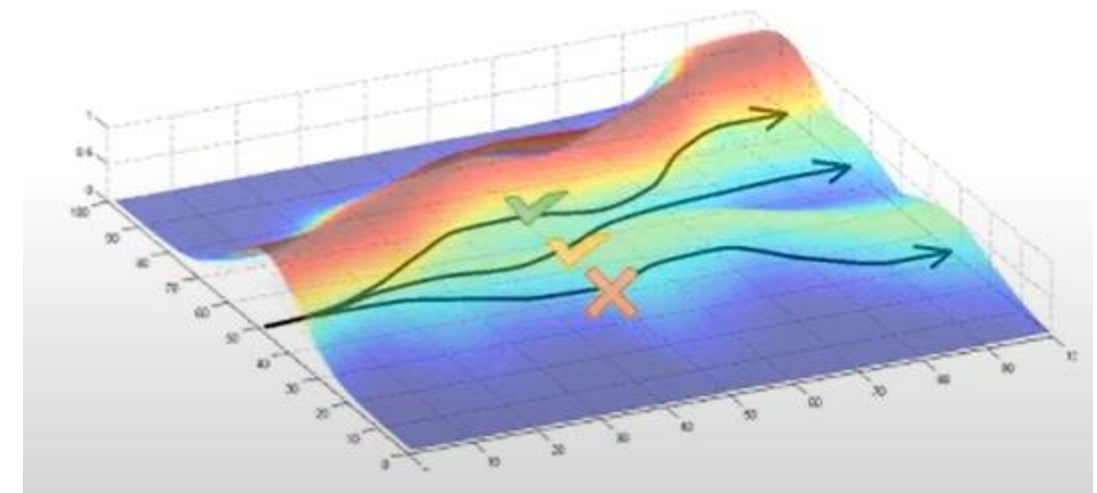
good stuff is made more likely

bad stuff is made less likely

simply formalizes the notion of “trial and error”!

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run it on the robot)
2. $\nabla_{\theta} J(\theta) \approx \sum_i \left(\sum_t \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_t^i | \mathbf{s}_t^i) \right) \left(\sum_t r(\mathbf{s}_t^i, \mathbf{a}_t^i) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$



REINFORCE SUGGESTED READING

- Classic papers
 - Williams (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning: introduces REINFORCE algorithm
 - Baxter & Bartlett (2001). Infinite-horizon policy-gradient estimation: temporally decomposed policy gradient (not the first paper on this! see actor-critic section later)
 - Peters & Schaal (2008). Reinforcement learning of motor skills with policy gradients: very accessible overview of optimal baselines and natural gradient
- Deep reinforcement learning policy gradient papers
 - Levine & Koltun (2013). Guided policy search: deep RL with importance sampled policy gradient (unrelated to later discussion of guided policy search)
 - Schulman, L., Moritz, Jordan, Abbeel (2015). Trust region policy optimization: deep RL with natural policy gradient and adaptive step size
 - Schulman, Wolski, Dhariwal, Radford, Klimov (2017). Proximal policy optimization algorithms: deep RL with importance sampled policy gradient

REINFORCE-IMPLEMENTATION

Algorithm 2.1 REINFORCE algorithm

```
1: Initialize learning rate  $\alpha$ 
2: Initialize weights  $\theta$  of a policy network  $\pi_\theta$ 
3: for  $episode = 0, \dots, MAX\_EPISODE$  do
4:   Sample a trajectory  $\tau = s_0, a_0, r_0, \dots, s_T, a_T, r_T$   $\longrightarrow$  Use the policy network  $\pi_\theta$  to generate a trajectory for an episode
5:   Set  $\nabla_\theta J(\pi_\theta) = 0$ 
6:   for  $t = 0, \dots, T$  do
7:      $R_t(\tau) = \sum_{t'=t}^T \gamma^{t'-t} r'_{t'}$   $\longrightarrow$  For each time step 't' in the trajectory, compute the return  $R_t(\tau)$ 
8:      $\nabla_\theta J(\pi_\theta) = \nabla_\theta J(\pi_\theta) + R_t(\tau) \nabla_\theta \log \pi_\theta(a_t | s_t)$   $\searrow$  Use  $R_t(\tau)$  to estimate the policy gradient and sum policy gradients for all time steps
9:   end for
10:   $\theta = \theta + \alpha \nabla_\theta J(\pi_\theta)$   $\longrightarrow$  Update the policy network parameters  $\theta$ 
11: end for
```

Note: A trajectory is discarded after each parameter update – it cannot be reused
 \rightarrow REINFORCE is an on-policy algorithm

REINFORCE-IMPLEMENTATION

❖ REINFORCE Training Loop

Code 2.6 REINFORCE implementation: training method

```
1  # slm_lab/agent/algorithm/reinforce.py
2
3  class Reinforce(Algorithm):
4      ...
5
6      @lab_api
7      def sample(self):
8          batch = self.body.memory.sample()
9          batch = util.to_torch_batch(batch, self.net.device,
10                                     ↪ self.body.memory.is_episodic)
11          return batch
12
13      @lab_api
14      def train(self):
15          ...
16          clock = self.body.env.clock
17          if self.to_train == 1:
18              batch = self.sample()
19              ...
20              pdparams = self.calc_pdparam_batch(batch)
21              advs = self.calc_ret_advs(batch)
22              loss = self.calc_policy_loss(batch, pdparams, advs)
23              self.net.train_step(loss, self.optim, self.lr_scheduler,
24                                  ↪ clock=clock, global_net=self.global_net)
25              # reset
26              self.to_train = 0
27              return loss.item()
28          else:
29              return np.nan
```

This code makes one parameter update to the policy network using a batch of collected trajectories

Computing the returns used to calculate the policy loss

The trajectories are obtained from the agent's memory by calling sample

Computing the action probability distribution

Updating the policy network parameters using the loss

REINFORCE-IMPLEMENTATION

❖ On-Policy Replay Memory

- REINFORCE is an on-policy algorithm
- Trajectories sampled by the algorithm should be discarded in a Memory class for learning, then deleted after each training step
- For implementation, the following API methods are contained
 1. `reset` clears and resets the memory class variables.
 2. `update` adds an experience to the memory.
 3. `sample` samples a batch of data for training.

REINFORCE-IMPLEMENTATION

❖ On-Policy Replay Memory : Memory Initialization and Reset

Code 2.7 OnPolicyReplay: reset

```
1 # slm_lab/agent/memory/onpolicy.py
2
3 class OnPolicyReplay(Memory):
4     ...
5
6     def __init__(self, memory_spec, body):
7         super().__init__(memory_spec, body)
8         # NOTE for OnPolicy replay, frequency = episode; for other classes
9         ↪ below frequency = frames
10        util.set_attr(self, self.body.agent.agent_spec['algorithm'],
11        ↪ ['training_frequency'])
12        # Don't want total experiences reset when memory is
13        self.is_episodic = True
14        self.size = 0 # total experiences stored
15        self.seen_size = 0 # total experiences seen cumulatively
16        # declare what data keys to store
17        self.data_keys = ['states', 'actions', 'rewards', 'next_states',
18        ↪ 'done']
19        self.reset()
20
21 @lab_api
22 def reset(self):
23     '''Resets the memory. Also used to initialize memory vars'''
24     for k in self.data_keys:
25         setattr(self, k, [])
26     self.cur_epi_data = {k: [] for k in self.data_keys}
27     self.most_recent = (None,) * len(self.data_keys)
28     self.size = 0
```

- “reset” in Code 2.7 is used to clear the memory after each training step
- This is specific to an on-policy memory because trajectories cannot be reused for later training

Individual episodes are constructed by storing experiences in the current episode data dictionary “self.cur_epi_data”

REINFORCE-IMPLEMENTATION

❖ On-Policy Replay Memory : Memory Update

Code 2.8 OnPolicyReplay: add experience

```
1  # slm_lab/agent/memory/onpolicy.py
2
3  class OnPolicyReplay(Memory):
4      ...
5
6      @lab_api
7      def update(self, state, action, reward, next_state, done):
8          '''Interface method to update memory'''
9          self.add_experience(state, action, reward, next_state, done)
10
11     def add_experience(self, state, action, reward, next_state, done):
12         '''Interface helper method for update() to add experience to memory'''
13         self.most_recent = (state, action, reward, next_state, done)
14         for idx, k in enumerate(self.data_keys):
15             self.cur_epi_data[k].append(self.most_recent[idx])
16         # If episode ended, add to memory and clear cur_epi_data
17         if util.epi_done(done):
18             for k in self.data_keys:
19                 getattr(self, k).append(self.cur_epi_data[k])
20             self.cur_epi_data = {k: [] for k in self.data_keys}
21             # If agent has collected the desired number of episodes, it is
22             #   ↳ ready to train
23             # length is num of epis due to nested structure
24             if len(self.states) ==
25                 ↳ self.body.agent.algorithm.training_frequency:
26                 self.body.agent.algorithm.to_train = 1
27
28     # Track memory size and num experiences
29     self.size += 1
30     self.seen_size += 1
```

Add the experience to the current episode

Check if the episode is finished. This is given by the "done" variable which will be '1' if the episode is finished, '0' otherwise

If the episode is finished, add the entire set of experiences for the episode to the main containers in the memory class

If the episode is finished, clear the current episode dictionary so that memory class is ready to store the next episode

If the desired number of episodes has been collected, set the 'train flag' to '1' in the agent. This signals that the agent should train this time step

REINFORCE-IMPLEMENTATION

❖ On-Policy Replay Memory : Memory Sample

Code 2.9 OnPolicyReplay: sample

```
1 # slm_lab/agent/memory/onpolicy.py
2
3 class OnPolicyReplay(Memory):
4
5     def sample(self):
6         batch = {k: getattr(self, k) for k in self.data_keys}
7         self.reset()
8         return batch
```

Returning all of the completed episodes packaged into a batch dictionary

Reset the memory as the stored experiences will no longer be valid once the agent has completed a training step

Assignment-Implementation of REINFORCE

- ❖ Implementing Code 2.1 and Code 2.10 in the textbook
 - Analyzing and interpreting codes line by line
 - Draw graph in Figure 2.2 (textbook) and explaining the results
- ❖ Reproducing the results in Code 2.12, Code 2.13, and Figure 2.3
 - Draw graphs in Figure 2.3 (textbook) and discussing the results
- ❖ Reproducing the results in Code 2.14, Code 2.15, and Figure 2.4
 - Draw graphs in Figure 2.4 (textbook) and discussing the results

Assignment-Implementation of REINFORCE

- ❖ Please submit the report (word format)
- ❖ Please submit all the codes from your implementation including
 - Main training codes
 - Reproduction codes of experimental results
 - Execution codes