



Additional Questions

Task 1.7 (1 pt): The code above for REINFORCE does not implement an explicit exploration strategy. Explain how the exploration-exploitation trade-off was circumvented by your choice of parameters initialisation in a few sentences, and how this approach could fail for some initialisations.

Write in the Markdown cell below (use LaTeX-math mode for equations, etc).

The exploration-exploitation trade-off in the REINFORCE implementation was implicitly managed through the initialization of the policy parameters θ . By initializing θ with small random values (e.g., uniformly between -0.1 and 0.1), the initial policy generated by the softmax function is approximately uniform across all actions in each state. This means that the agent starts with a high degree of exploration, as each action has nearly equal probability of being selected. Over time, as the agent collects experience and updates θ based on received rewards, the policy naturally shifts towards exploiting actions that yield higher returns.

This approach can fail with certain initializations. If θ is initialized with large magnitudes or biased values, the softmax function may produce a policy that heavily favors certain actions from the start, reducing exploration. For instance, if one action's parameter is significantly higher than others due to initialization, the agent might consistently select that action and miss out on exploring potentially better alternatives. Therefore, without an explicit exploration strategy, the agent's ability to explore effectively depends on careful initialization of θ , and inappropriate initial values could hinder learning by limiting exploration.

Task 1.8 (0.5 pt): In a few sentences explain the shortcomings of the direct parametrisation.

Write in the Markdown cell below (use LaTeX-math mode for equations, etc).

The direct parametrization has significant shortcomings due to its lack of scalability and generalization. By assigning a unique parameter to each state-action pair, it results in a parameter vector whose size grows exponentially with the number of states and actions. This makes it impractical for environments with large or continuous state spaces. Additionally, it doesn't capture similarities or shared structures between different states or actions, preventing the agent from generalizing learned behaviors to unseen situations. Consequently, learning becomes inefficient, requiring extensive data to cover all possible state-action combinations.

Part 2: CartPole with A2C

The Cart Pole Environment

The **CartPole** environment is a classical problem where a pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. The pendulum is placed upright on the cart and the goal is to balance the pole by applying forces in the left and right direction on the cart.

Since the focus of this task is the implementation of reinforcement learning algorithms, it is not necessary to have a detailed understanding of the mechanics of the Cart Pole environment beyond the observation and action space sizes.

Initialize the environment

Initialize the Cart Pole environment in the Gymnasium library to start learning a policy for it. You can read more about the Gymnasium library [here](#).

```
In [20]: env = gym.make("CartPole-v1")
```

PyTorch Intro

We also introduce basic functionality for PyTorch. PyTorch is a widely used automatic differentiation library, which is very useful for training neural networks. This part is also for guidance and does not include questions. While we do introduce the basics of PyTorch, you are encouraged to explore more about the library following the official [tutorials](#) in case you haven't used it before.

PyTorch is built on top of the class `torch.Tensor`, which implements many operations on vectors, matrices and higher dimensional tensors. The functionality of operations on torch tensors closely mirrors the operations from the `numpy` package.

```
In [21]: torch.zeros(5)
```

```
Out[21]: tensor([0., 0., 0., 0., 0.])
```

```
In [22]: a_matrix = [[1., 2.], [3., 4.]]
          torch.tensor(a_matrix)
```

```
Out[22]: tensor([[1., 2.],
                [3., 4.]])
```

```
In [23]: torch.matmul(torch.tensor(a_matrix),
                       torch.tensor([4., 5.]))
```

```
Out[23]: tensor([14., 32.])
```

What makes torch special is that it implements automated differentiation by tracking the operation history on tensors. Namely, if a `torch.Tensor` has the property `requires_grad` set to `True`, `torch` will automatically track operations made on that tensor. A final `.backward()` call to a computation result will compute partial derivatives of the variable with respect to all tensors in the tracked computational graph. This is particularly useful for tracking derivatives with respect to a loss function, to enable simple gradient descent.

```
In [24]: # example from https://pytorch.org/tutorials/beginner/basics/autogradqs_tutorial.html

x = torch.ones(5) # input tensor
y = torch.zeros(3) # expected output
w = torch.randn(5, 3, requires_grad=True) # weights
b = torch.randn(3, requires_grad=True) # bias
z = torch.matmul(x, w)+b # Linear combination
loss = torch.nn.functional.binary_cross_entropy_with_logits(z, y)
print(loss)
```

```
tensor(1.0208, grad_fn=<BinaryCrossEntropyWithLogitsBackward0>)
```

```
In [25]: # differentiate with respect to loss
          loss.backward()
          # the computed partial derivatives of w
```

```
Out[25]: tensor([0.1598, 0.2628, 0.1917],
                [0.1598, 0.2628, 0.1917],
                [0.1598, 0.2628, 0.1917],
                [0.1598, 0.2628, 0.1917],
                [0.1598, 0.2628, 0.1917])
```

To prevent tracking history, you can also wrap a piece of code in `with torch.no_grad():`. This can be particularly helpful when evaluating a model on new data, when training (and hence differentiation) is not required.

For more advanced modeling (as will be required by this forward), it is usually useful to implement classes as children of `nn.Module`. These implement the function `forward`, which returns the output of the forward pass of a model given inputs. Modules have `nn.Parameter`, which will be the trainable parameters of the models. Most useful sub-modules are already provided by the torch package: see for instance `nn.Linear`.

Finally, torch implements many standard optimizers: for instance `torch.optim.SGD`, `torch.optim.Adam`. These are typically used to automatically take care of updating some partial derivatives are computed.

Advantage Actor-Critic (A2C)

A more stable alternative to vanilla policy gradients via REINFORCE is the so-called A2C algorithm, which is a actor-critic method. First, you should familiarise yourself with the A2C algorithm reading the original [paper](#).

Task 2.1 (3 pts): Explain in a few sentences each,

1. how the A2C algorithm addresses the shortcomings of vanilla policy gradient methods,
2. the difference between A2C and A3C,
3. why the A3C algorithm could be expected to perform even better than A2C.

Write in the Markdown cell below (use LaTeX-math mode for equations, etc).

1. How the A2C algorithm addresses the shortcomings of vanilla policy gradient methods:

The A2C (Advantage Actor-Critic) algorithm improves upon vanilla policy gradient methods like REINFORCE by reducing variance and enhancing learning stability. In vanilla policy gradients, the policy is updated using returns (the cumulative future rewards), which can have high variance and lead to unstable learning. A2C introduces a critic network that estimates the value function $V(s)$, representing the expected return from state s . By using the advantage function $A(s,a) = Q(s,a) - V(s)$, where $Q(s,a)$ is the expected return after taking action a in state s , the algorithm focuses updates on actions that perform better than expected. This approach reduces the variance in gradient estimates by subtracting the baseline value $V(s)$, leading to more stable and efficient learning.

1. The difference between A2C and A3C:

A2C and A3C (Asynchronous Advantage Actor-Critic) are both actor-critic algorithms that utilize the advantage function to update policies. The primary difference lies in how they handle parallelism and updates:

A2C: A2C is a synchronous version of A3C. It runs multiple worker processes in parallel environments, but synchronizes the gradient updates. After a set number of steps, all workers send their gradients to a central optimizer, which averages them and updates the global network parameters synchronously. This method benefits from efficient utilization of GPUs and simplifies implementation compared to asynchronous methods.

A3C: A3C operates with multiple workers running asynchronously. Each worker interacts with its own environment and updates the global network parameters independently without waiting for other workers. The asynchronous updates introduce diversity in learning and help prevent workers from becoming stuck in the same local optima. This approach leverages CPU resources effectively and can lead to faster learning due to increased exploration.

1. Why the A3C algorithm could be expected to perform even better than A2C:

A3C can potentially outperform A2C because its asynchronous nature introduces greater exploration and diversity in learning. As workers operate independently, they experience different state transitions and update the global parameters based on varied experiences. This diversity helps in:

Escaping Local Optima: Independent exploration increases the chances of discovering better policies that might be missed in synchronous updates.

Reducing Parameter Update Conflicts: Asynchronous updates minimize the risk of workers overwriting each other's progress in a synchronized manner.

Improving Learning Efficiency: Without the need to wait for synchronization, workers can perform updates more frequently, leading to faster convergence.

Additionally, the inherent randomness and non-determinism in asynchronous updates can act as a form of regularization, potentially improving the generalization of the learned policy.

Neural Network Parameterization

We will use neural networks to parameterize both the value function and the policy. The **ActorNet** defines the policy $\pi_{\theta}(a|s)$. Given state s as input, it will output a distribution over the action space. The **CriticNet** defines the approximation to the value function, which could be denoted $\hat{V}_{\omega}(s)$. Given state s as input, it will output an approximation of the state-value.

You can play with different network architectures and the **activation** functions. Note that the **ActorNet** should produce a probability distribution over the action space.

Task 2.2 (0.5 pt): Design your own neural networks for **ActorNet** and **CriticNet** in the code block below.

```
In [26]: class ActorNet(torch.nn.Module):
          def __init__(self, num_state: int, num_action: int):
              super(ActorNet, self).__init__()

              # Define the network architecture
              self.model = torch.nn.Sequential(
                  torch.nn.Linear(num_state, 128), # Input layer to hidden layer
                  torch.nn.ReLU(), # Activation function
                  torch.nn.Linear(128, num_action), # Hidden layer to output layer
                  torch.nn.Softmax(dim=-1) # Softmax to get probability distribution
              )

          # Return a probability distribution over the action space
          def forward(self, state):
              return self.model(state)

          class CriticNet(torch.nn.Module):
              def __init__(self, num_state: int):
                  super(CriticNet, self).__init__()

              # Define the network architecture
              self.model = torch.nn.Sequential(
                  torch.nn.Linear(num_state, 128), # Input layer to hidden layer
                  torch.nn.ReLU(), # Activation function
                  torch.nn.Linear(128, 1) # Hidden layer to output layer (single value)
              )

              # Return a single value
              def forward(self, state):
                  return self.model(state)
```

Learning Actor and Critic

We train both the actor and critic nets by interacting with the environment. You implement the A2C training process below.

Task 2.3 (4 pts): Fill in and run the code below to implement A2C in the Cart Pole environment. You should play with your hyper-parameters as well, to enable efficient learning. You may play with different learning rates and optimizers.

For different optimizers in PyTorch, see [here](#).

```
In [35]: import torch.nn as nn
          import torch.nn.functional as F

          class A2C():

              def __init__(self, env):
                  """Initialization code"""
                  # Initiate both nets
                  self.actor = ActorNet(env.observation_space.shape[0], env.action_space.n)
                  self.critic = CriticNet(env.observation_space.shape[0])

                  # Set gamma and learning rates
                  self.gamma = 0.99 # Discount factor
                  actor_lr = 0.005 # Learning rate for the actor
                  critic_lr = 0.005 # Learning rate for the critic
                  self_entropy_coef = 0.0001 # Entropy regularization coefficient (tau)
                  self.env = env

                  # You can also experiment with different optimizers
                  self.actor_opt = torch.optim.Adam(self.actor.parameters(), lr=actor_lr)
                  self.critic_opt = torch.optim.Adam(self.critic.parameters(), lr=critic_lr)

              def train(self, max_iter: int = 1000, max_episode_len: int = 500):

                  total_reward = []

                  pbar = tqdm.tqdm(range(max_iter), desc="Episode")
                  for num_iter in pbar:
                      # At each iteration, we roll out an episode using current policy
                      rewards = []
                      states = []
                      actions = []
                      log_probs = []
                      values = []

                      s, _ = self.env.reset()
                      for t in range(max_episode_len):
                          # Given current state, get action probabilities from actor
                          state = torch.tensor(s, dtype=torch.float32).unsqueeze(0)
                          probs = self.actor(state)
                          # Create the distribution as current policy
                          pi = torch.distributions.Categorical(probs)
                          # Sample one action from this policy
                          a = pi.sample()

                          log_prob = pi.log_prob(a)
                          value = self.critic(state)

                          # Interact with the environment
                          s, r, is_terminal, _ = self.env.step(a.item())

                          # Store data
                          rewards.append(r)
                          states.append(state)
                          actions.append(a)
                          log_probs.append(log_prob)
                          values.append(value)

                          if is_terminal:
                              break

                      # Store total_reward
                      total_reward.append(sum(rewards))

                      pbar.set_postfix_str(f'Last reward: {total_reward[-1]}, Mean last (min(last_reward, total_reward))')

                      # Compute returns and advantages
                      returns = []
                      R = 0
                      for r in reversed(rewards):
                          R = r + self.gamma * R
                          returns.insert(0, R)
                      returns = torch.tensor(returns, dtype=torch.float32)

                      # Convert lists to tensors
                      log_probs = torch.stack(log_probs)
                      values = torch.stack(values).squeeze()

                      # Compute advantages
                      advantages = returns - values.detach()

                      # Compute actor loss (policy gradient with baseline)
                      actor_loss = - (log_probs * advantages).mean()

                      # Compute critic loss (value function approximation)
                      critic_loss = F.mse_loss(values, returns)

                      # Update actor network
                      self.actor_opt.zero_grad()
                      actor_loss.backward()
                      self.actor_opt.step()

                      # Update critic network
                      self.critic_opt.zero_grad()
                      critic_loss.backward()
                      self.critic_opt.step()

                      self.env.close()

                  return total_reward
```

```
In [36]: # TODO: Run your implementation here.
          max_iter = 1000
          max_episode_len = 500

          a2c = A2C(env)

          total_reward = a2c.train(max_iter, max_episode_len)

          Episode: 100% | 1000/1000 | 05:43:00:00, 2.91it/s, Last reward: 500.0, Mean last 100 episodes: 499.63
```

Useful hints:

1. See [here](#) for distributions in PyTorch. You may find the function `log_prob()` useful. You can easily compute $\log \pi(a|s)$ and its gradient using the function.
2. To optimize a model in PyTorch, you first define the loss function, then call the corresponding optimizer and perform backpropagation. For example, to optimize over the value function approximation, first let $\text{loss} = \|V_{\omega} - V^*\|^2$ and then call `loss.backward()` to compute the gradient.
3. You may find the function `detach()` defined in PyTorch useful if you need to call `backward()` multiple times on the same variables without storing the computational graph every time.
4. Maximizing $f(x)$ is equivalent to minimizing $-f(x)$. The default optimizers in PyTorch perform minimization.
5. You may find the classes `torch.nn.Softmax`, `torch.nn.ReLU`, `torch.nn.Tanh` useful when implementing a neural network architecture.

Analyze your results

You can run the codes below to test your algorithm. These will help you understand if your algorithm is working, and will help us grade your implementation.

```
In [31]: # This will visualize your learnt policy on the Cart Pole environment
          from IPython.display import clear_output

          env_render = gym.make("CartPole-v1", render_mode="rgb_array")
          for num_episode in range(1):
              s, _ = env_render.reset()
              tot_reward = []
              for t in range(max_episode_len):
                  state = torch.tensor(s, dtype=torch.float).view(1, -1)

                  prob = a2c.actor(state)
                  policy = torch.distributions.categorical.Categorical(prob)
                  action = policy.sample()
                  s, r, is_terminal, _ = env_render.step(action.item())
                  tot_reward.append(r)
                  if is_terminal:
                      print(f"Episode {t} terminated after {t} steps with reward {r}. ", format(num_episode, 1))
                      break
                  clear_output(wait=True)
                  plt.imshow(env_render.render())
                  plt.show()

          env_render.close()
```


Episode 0 terminated after 318 steps with reward 318.0.

We use a simple criteria to judge whether you solve the problem. We will compute an average of the total reward over the previous 100 episodes at each iteration. If there exists an iteration with average total reward larger than a threshold, as shown below, we will give you full grades. That is, the problem is solved when the average reward is greater than or equal to 475 over 100 consecutive trials.

Note that **CartPole-v1** has a termination condition of 500 timesteps. It's done so that one episode doesn't take forever. So we say that if a policy can balance a pole for 500 time steps (and achieve 500 reward) it's probably good enough.

```
In [37]: # this code should print success
          success_flag = False
          buffer_len = 100
          buffer_sum = sum(total_reward[:buffer_len])
          avg_reward = [buffer_sum / buffer_len]

          for num_iter in range(buffer_len, max_iter):
              buffer_sum += total_reward[num_iter]
              buffer_sum -= total_reward[num_iter - buffer_len]
              buffer_avg = buffer_sum / buffer_len
              avg_reward.append(buffer_avg)

              if buffer_avg >= env.spec.reward_threshold and not success_flag:
                  print(f'Successfully solve the problem in {t} iterations.', format(num_iter + 1))
                  success_flag = True

          if not success_flag:
              print('Unfortunately the agent is not smart enough.')
```

```
# plot the average reward curve through the training process
plt.plot(range(buffer_len-1, len(total_reward)), avg_reward)
plt.show()
```


Successfully solve the problem in 422 iterations.

```
In [38]: class A2CEntropy():
          def __init__(self, env):
              """Initialization code"""
              # Initiate both nets
              self.actor = ActorNet(env.observation_space.shape[0], env.action_space.n)
              self.critic = CriticNet(env.observation_space.shape[0])

              # Set gamma and learning rates
              self.gamma = 0.99 # Discount factor
              actor_lr = 0.005 # Learning rate for the actor
              critic_lr = 0.005 # Learning rate for the critic
              self_entropy_coef = 0.0001 # Entropy regularization coefficient (tau)
              self.env = env

              # You can also experiment with different optimizers
              self.actor_opt = torch.optim.Adam(self.actor.parameters(), lr=actor_lr)
              self.critic_opt = torch.optim.Adam(self.critic.parameters(), lr=critic_lr)

              self_entropy_coef = 0.0001 # Entropy regularization coefficient (tau)

              def train(self, max_iter: int = 1000, max_episode_len: int = 500):

                  total_reward = []

                  pbar = tqdm.tqdm(range(max_iter), desc="Episode")
                  for num_iter in pbar:
                      # At each iteration, we roll out an episode using current policy
                      rewards = []
                      log_probs = []
                      values = []
                      entropies = []

                      s, _ = self.env.reset()
                      for t in range(max_episode_len):
                          # Given current state, get action probabilities from actor
                          state = torch.tensor(s, dtype=torch.float32).unsqueeze(0)
                          probs = self.actor(state)
                          # Create the distribution as current policy
                          pi = torch.distributions.Categorical(probs)
                          # Sample one action from this policy
                          a = pi.sample()

                          log_prob = pi.log_prob(a)
                          entropy = pi.entropy()

                          value = self.critic(state)

                          # Interact with the environment
                          s_new, r, is_terminal, _ = self.env.step(a.item())

                          # Store data
                          rewards.append(r)
                          log_probs.append(log_prob)
                          values.append(value)
                          entropies.append(entropy)

                          if is_terminal:
                              break

                      # Store total_reward
                      total_reward.append(sum(rewards))

                      mean_reward = sum(total_reward[-100:]) / len(total_reward)

                      # Convert lists to tensors
                      log_probs = torch.stack(log_probs)
                      values = torch.stack(values).squeeze()
                      entropies = torch.stack(entropies) # Convert entropies to tensor

                      # Now compute average entropy
                      avg_entropy = entropies.mean().item()

                      pbar.set_postfix_str(f'Last reward: {total_reward[-1]}, Avg Entropy: {avg_entropy}')

                      # Compute returns and advantages
                      returns = []
                      R = 0
                      for r in reversed(rewards):
                          R = r + self.gamma * R
                          returns.insert(0, R)
                      returns = torch.tensor(returns, dtype=torch.float32)

                      # Compute advantages
                      advantages = returns - values.detach()

                      # Compute actor loss (policy gradient with baseline and entropy regularization)
                      actor_loss = - (log_probs * advantages).mean() + self_entropy_coef * avg_entropy

                      # Compute critic loss (value function approximation)
                      critic_loss = F.mse_loss(values, returns)

                      # Update actor network
                      self.actor_opt.zero_grad()
                      actor_loss.backward()
                      self.actor_opt.step()

                      # Update critic network
                      self.critic_opt.zero_grad()
                      critic_loss.backward()
                      self.critic_opt.step()

                      self.env.close()

                  return total_reward
```

```
In [39]: # TODO: Run your implementation here.
          max_iter_entropy = 1000
          max_episode_len_entropy = 500

          a2c_entropy = A2CEntropy(env)

          total_reward_entropy = a2c_entropy.train(max_iter_entropy, max_episode_len_entropy)

          Episode: 100% | 1000/1000 | 06:14:00:00, 2.67it/s, Last reward: 500.0, Avg Entropy: 0.4996, Mean last 100 episodes: 50
```

```
In [40]: # this code should print success
          success_flag = False
          buffer_len = 100
          buffer_sum = sum(total_reward_entropy[:buffer_len])
          avg_reward = [buffer_sum / buffer_len]

          for num_iter in range(buffer_len, max_iter):
              buffer_sum += total_reward_entropy[num_iter]
              buffer_sum -= total_reward_entropy[num_iter - buffer_len]
              buffer_avg = buffer_sum / buffer_len
              avg_reward_entropy.append(buffer_avg)

              if buffer_avg >= env.spec.reward_threshold and not success_flag:
                  print(f'Successfully solve the problem in {t} iterations.', format(num_iter + 1))
                  success_flag = True

          if not success_flag:
              print('Unfortunately the agent is not smart enough.')
```

```
# plot the average reward curve through the training process
plt.plot(range(buffer_len-1, len(total_reward_entropy)), avg_reward)
plt.show()
```


Successfully solve the problem in 497 iterations.

```
In [41]: # TODO: Visualize your results and compare it to the results without entropy in Task 2.3
          import numpy as np

          def moving_average(data, window_size):
              return np.convolve(data, np.ones(window_size)/window_size, mode='valid')

          # Define the window size for the moving average
          window_size = 30

          # Compute the moving averages
          smoothed_total_reward = moving_average(total_reward, window_size)
          smoothed_total_reward_entropy = moving_average(total_reward_entropy, window_size)

          # Adjust the x-axis to match the length of the smoothed data
          episodes = np.arange(window_size - 1, len(total_reward))

          # Plotting the smoothed results
          plt.figure(figsize=(12, 6))

          # Plot smoothed total rewards for the A2C agent without entropy regularization
          plt.plot(episodes, smoothed_total_reward, label='A2C without Entropy Regularization')

          # Plot smoothed total rewards for the A2C agent with entropy regularization
          plt.plot(episodes, smoothed_total_reward_entropy, label='A2C with Entropy Regularization')

          # Adding labels and title
          plt.xlabel('Episode')
          plt.ylabel('Total Reward (Smoothed)')
          plt.title('Smoothed Comparison of A2C Training Performance on CartPole-v1')

          # Adding legend
          plt.legend()

          # Display the plot
          plt.grid(True)
          plt.show()
```


Successfully solve the problem in 497 iterations.