

Deep Reinforcement Learning

Shan-Hung Wu & DataLab
Fall 2025

In the last lab, we use the tabular method (Q-learning, SARSA) to train an agent to play *Flappy Bird* with features in environments. However, it is time-costly and inefficient if more features are added to the environment because the agent can not easily generalize its experience to other states that were not seen before. Furthermore, in realistic environments with large state / action space, it requires a large memory space to store all state-action pairs.

In this lab, we introduce deep reinforcement learning, which utilizes function approximation to estimate value / policy for all unseen states such that given a state, we can estimate its value or action. We can use what we have learned in machine learning (e.g. regression, DNN) to achieve it.

PPO X GAE

Reference: Generalized Advantage Estimation, Proximal Policy Optimization

To use reinforcement learning successfully in situations approaching real-world complexity, however, agents are confronted with a difficult task: they must derive efficient representations of the environment from high-dimensional sensory inputs, and use these to generalize past experience to new situations.

In this lab, we are going to train an agent which takes raw frames as input instead of hand-crafted features.

```
In [ ]: import numpy as np
import matplotlib.pyplot as plt
import moviepy.editor as mpy
import skimage.transform
from IPython.display import Image, display

import tensorflow as tf
import tensorflow_probability as tfp
import tensorflow.keras.losses as kls
```

```
In [ ]: gpus = tf.config.list_physical_devices("GPU")
if gpus:
    try:
        # Restrict TensorFlow to only use the fourth GPU
        tf.config.set_visible_devices(gpus[0], 'GPU')

        # Currently, memory growth needs to be the same across GPUs
        for gpu in gpus:
            tf.config.experimental.set_memory_growth(gpu, True)

    logical_gpus = tf.config.list_logical_devices('GPU')
```

```

        print(len(gpus), "Physical GPUs,", len(logical_gpus), "Logical GPUs")
    except RuntimeError as e:
        # Memory growth must be set before GPUs have been initialized
        print(e)

```

```

In [ ]: import os
os.environ["SDL_VIDEODRIVER"] = "dummy" # this line make pop-out window not appear
from ple.games.flappybird import FlappyBird
from ple import PLE

game = FlappyBird()
env = PLE(game, fps=30, display_screen=False) # environment interface to game
env.reset_game()

test_game = FlappyBird()
test_env = PLE(test_game, fps=30, display_screen=False)
test_env.reset_game()

```

```

In [ ]: path = './movie_f'
if not os.path.exists(path):
    os.makedirs(path)

```

```

In [ ]: hparas = {
    'image_size': 84,
    'num_stack': 4,
    'action_dim': len(env.getActionSet()),
    'hidden_size': 256,
    'lr': 0.0001,
    'gamma': 0.99,
    'lambda': 0.95,
    'clip_val': 0.2,
    'ppo_epochs': 8,
    'test_epochs': 1,
    'num_steps': 512,
    'mini_batch_size': 64,
    'target_reward': 200,
    'max_episode': 30000,
}

```

```

In [ ]: # Please do not modify this method
def make_anim(images, fps=60, true_image=False):
    duration = len(images) / fps

    def make_frame(t):
        try:
            x = images[int(len(images) / duration * t)]
        except:
            x = images[-1]

        if true_image:
            return x.astype(np.uint8)
        else:
            return ((x + 1) / 2 * 255).astype(np.uint8)

    clip = mpy.VideoClip(make_frame, duration=duration)
    clip.fps = fps

    return clip

```

```
In [ ]: def preprocess_screen(screen):
    screen = skimage.transform.rotate(screen, -90, resize=True)
    screen = screen[:400, :]
    screen = skimage.transform.resize(screen, [hparas['image_size'], hparas['image_size']])
    return screen.astype(np.float32)

def frames_to_state(input_frames):
    if(len(input_frames) == 1):
        state = np.concatenate(input_frames*4, axis=-1)
    elif(len(input_frames) == 2):
        state = np.concatenate(input_frames[0:1]*2 + input_frames[1:]*2, axis=-1)
    elif(len(input_frames) == 3):
        state = np.concatenate(input_frames + input_frames[2:], axis=-1)
    else:
        state = np.concatenate(input_frames[-4:], axis=-1)

    return state
```

```
In [ ]: class ActorCriticNetwork(tf.keras.Model):
    def __init__(self, hparas):
        super().__init__()

        self.feature_extractor = tf.keras.Sequential([
            # Convolutional Layers
            tf.keras.layers.Conv2D(filters=32, kernel_size=8, strides=4),
            tf.keras.layers.ReLU(),
            tf.keras.layers.Conv2D(filters=64, kernel_size=4, strides=2),
            tf.keras.layers.ReLU(),
            tf.keras.layers.Conv2D(filters=64, kernel_size=3, strides=1),
            tf.keras.layers.ReLU(),
            # Embedding Layers
            tf.keras.layers.Flatten(),
            tf.keras.layers.Dense(hparas['hidden_size']),
            tf.keras.layers.ReLU(),
        ])

        # Actor Network
        self.actor = tf.keras.layers.Dense(hparas['action_dim'], activation='softmax')
        # Critic Network
        self.critic = tf.keras.layers.Dense(1, activation = None)

    def call(self, input):
        x = self.feature_extractor(input)
        action_logits = self.actor(x)
        value = self.critic(x)
        return action_logits, value
```

```
In [ ]: class Agent():
    def __init__(self, hparas):
        self.gamma = hparas['gamma']
        self.optimizer = tf.keras.optimizers.Adam(learning_rate=hparas['lr'])
        self.actor_critic = ActorCriticNetwork(hparas)
        self.clip_pram = hparas['clip_val']

    def ppo_iter(self, mini_batch_size, states, actions, log_probs, returns, adv_batch_size = states.shape[0]):
        for _ in range(batch_size // mini_batch_size):
            rand_ids = tf.convert_to_tensor(np.random.randint(0, batch_size, mini_batch_size))
            yield tf.gather(states, rand_ids), tf.gather(actions, rand_ids), tf.gather(log_probs, rand_ids), tf.gather(returns, rand_ids), tf.gather(advantages, rand_ids)
```

```

        tf.gather(returns, rand_ids), tf.gather(advantage, rand_ids)

def ppo_update(self, ppo_epochs, mini_batch_size, states, actions, log_probs
    total_actor_loss = 0
    total_critic_loss = 0
    for _ in range(ppo_epochs):
        for state, action, old_log_probs, reward, advantage in self.ppo_iter
            reward = tf.expand_dims(reward, axis=-1)

            with tf.GradientTape() as tape:
                prob, value = self.actor_critic(state, training=True)
                dist = tfp.distributions.Categorical(probs=prob, dtype=tf.fl
                entropy = tf.math.reduce_mean(dist.entropy())
                new_log_probs = dist.log_prob(action)

            # PPO ratio
            ratio = tf.math.exp(new_log_probs - old_log_probs)
            surr1 = ratio * advantage
            surr2 = tf.clip_by_value(ratio, 1.0 - self.clip_pram, 1.0 +

            actor_loss = tf.math.negative(tf.math.reduce_mean(tf.math.mi
            critic_loss = 0.5 * tf.math.reduce_mean(kls.mean_squared_err

            total_loss = actor_loss + critic_loss

            # single optimizer
            grads = tape.gradient(total_loss, self.actor_critic.trainable_va
            self.optimizer.apply_gradients(zip(grads, self.actor_critic.trai

            total_actor_loss += actor_loss
            total_critic_loss += critic_loss
    return total_actor_loss, total_critic_loss

```

```

In [ ]: # https://arxiv.org/pdf/1506.02438.pdf
# Equation 16
def compute_gae(rewards, masks, values, gamma, LAMBDA):
    gae = 0
    returns = []
    for i in reversed(range(len(rewards))):
        delta = rewards[i] + gamma * values[i + 1] * masks[i] - values[i]
        gae = delta + gamma * LAMBDA * masks[i] * gae
        returns.append(gae + values[i])

    returns.reverse()
    return returns

```

Testing Environment

```

In [ ]: def test_reward(test_env, agent):
    total_reward = 0
    # Reset the environment
    test_env.reset_game()
    input_frames = [preprocess_screen(test_env.getScreenGrayscale())]

    while not test_env.game_over():

        state = frames_to_state(input_frames)
        state = tf.expand_dims(state, axis=0)

```

```

        prob, value = agent.actor_critic(state)

        action = np.argmax(prob[0].numpy())
        reward = test_env.act(test_env.getActionSet()[action])
        total_reward += reward

        input_frames.append(preprocess_screen(test_env.getScreenGrayscale()))

    return total_reward

```

Training

```

In [ ]: agent = Agent(hparams)
max_episode = hparams['max_episode']
test_per_n_episode = 10
force_save_per_n_episode = 1000
early_stop_reward = 10

start_s = 0
best_reward = -5.0

checkpoint = tf.train.Checkpoint(
    actor_critic = agent.actor_critic,
    optimizer = agent.optimizer,
)

# Load from old checkpoint
# checkpoint.restore('ckpt_dir/ckpt-?')

```

```

In [ ]: ep_reward = []
total_avgr = []
early_stop = False
avg_rewards_list = []

env.reset_game()

for s in range(0, max_episode):
    if early_stop == True:
        break

    rewards = []
    states = []
    actions = []
    log_probs = []
    masks = []
    values = []

    display_frames = [env.getScreenRGB()]
    input_frames = [preprocess_screen(env.getScreenGrayscale())]

    for step in range(hparams['num_steps']):

        state = frames_to_state(input_frames)
        state = tf.expand_dims(state, axis=0)
        prob, value = agent.actor_critic(state)

        dist = tfp.distributions.Categorical(probs=prob[0], dtype=tf.float32)
        action = dist.sample(1)

```

```

        log_prob = dist.log_prob(action)

        reward = env.act(env.getActionSet()[int(action.numpy())])

        done = env.game_over()

        states.append(state)
        actions.append(action)
        values.append(value[0])
        log_probs.append(log_prob)
        rewards.append(tf.convert_to_tensor(reward, dtype=tf.float32))
        masks.append(tf.convert_to_tensor(1-int(done), dtype=tf.float32))

        display_frames.append(env.getScreenRGB())
        input_frames.append(preprocess_screen(env.getScreenGrayscale()))

    if done:
        env.reset_game()
        input_frames = [preprocess_screen(env.getScreenGrayscale())]

    _, next_value = agent.actor_critic(state)
    values.append(next_value[0])

    returns = compute_gae(rewards, masks, values, hparams['gamma'], hparams['lambda'])

    returns = tf.concat(returns, axis=0)
    log_probs = tf.concat(log_probs, axis=0)
    values = tf.concat(values, axis=0)
    states = tf.concat(states, axis=0)
    actions = tf.concat(actions, axis=0)
    advantage = returns - values[:-1]

    a_loss, c_loss = agent.ppo_update(hparams['ppo_epochs'], hparams['mini_batch_size'])
    print('[Episode %d] Actor loss: %.5f, Critic loss: %.5f' % (s, a_loss, c_loss))

    if s % test_per_n_episode == 0:
        # test agent hparams['test_epochs'] times to get the average reward
        avg_reward = np.mean([test_reward(test_env, agent) for _ in range(hparams['test_epochs'])])
        print("Test average reward is %.1f, Current best average reward is %.1f" % (avg_reward, best_reward))
        avg_rewards_list.append(avg_reward)

        if avg_reward > best_reward:
            best_reward = avg_reward
            agent.actor_critic.save('./save/Actor/model_actor_{0}_{1}'.format(s, avg_reward))
            checkpoint.save(file_prefix = './save/checkpoints/ckpt')

    if s % force_save_per_n_episode == 0:
        agent.actor_critic.save('./save/Actor/model_actor_{0}_{1}'.format(s, avg_reward))
        checkpoint.save(file_prefix = './save/checkpoints/ckpt')
        clip = make_anim(display_frames, fps=60, true_image=True).rotate(-90)
        clip.write_videofile("movie_f/{0}_demo-{1}.webm".format('Lab15', s), fps=60)
        display(clip.ipython_display(fps=60, autoplay=1, loop=1, maxduration=120))

    if best_reward >= early_stop_reward:
        early_stop = True

```

Assignment

What you should do:

- Run the code and comprehend it
- Write your discovery in this notebook (exempli gratia how many times your birds fly to get more than 10 rewards)

Evaluation metrics:

- Report of this lab (50%)
- The bird is able to fly through at least 1 pipe (50%)

Requirements:

- Upload the notebook to eclass
 - Lab15_{student_id}.ipynb
- **Deadline: 2025-12-10 (Wed) 23:59**