

cs285 hw1

- cs285 hw1
 - Code
 - Atari
 - DQN
 - Pytorch
 - RL framework design
 - Exps
 - DQN
 - AC

Code

Atari

1. DeepMind-style Atari env:

```
def wrap_deepmind(env):
    """Make end-of-life == end-of-episode, but only reset on true game over. Done by Dee
    since it helps value estimation. 没有命了立马重启."""
    env = EpisodicLifeEnv(env)
    """Sample initial states by taking random number of no-ops on reset. No-op is assume
    env = NoopResetEnv(env, noop_max=30)
    """Return only every `skip`-th frame"""
    env = MaxAndSkipEnv(env, skip=4)
    """Take action on reset for environments that are fixed until firing."""
    if 'FIRE' in env.unwrapped.get_action_meanings():
        env = FireResetEnv(env)
    """把原始图片处理成84*84*1的大小"""
    env = ProcessFrame84(env)
    """Bin reward to {+1, 0, -1} by its sign."""
    env = ClipRewardEnv(env)
    return env
```

2. config

```
{
    'learning_starts': 50000,
    'target_update_freq': 10000,
    'replay_buffer_size': int(1e6),
    'num_timesteps': int(2e8),
    'learning_freq': 4, # train freq, 和'frame_history_len'一致
    'grad_norm_clipping': 10,
    'frame_history_len': 4, # 合并相邻4 frames to represent state (84, 84, 4)
}
```

atari optimizer config:

```
def atari_optimizer(num_timesteps):
    lr_schedule = PiecewiseSchedule(
        [
            (0, 1e-1),
            (num_timesteps / 40, 1e-1),
            (num_timesteps / 8, 5e-2),
        ],
        outside_value=5e-2,
    )

    return OptimizerSpec(
        constructor=optim.Adam,
        optim_kwargs=dict(lr=1e-3, eps=1e-4),
        learning_rate_schedule=lambda t: lr_schedule.value(t),
    )
```

atari explore config:

```
def atari_exploration_schedule(num_timesteps):
    return PiecewiseSchedule([
        (0, 1.0),
        (1e6, 0.1),
        (num_timesteps / 8, 0.01),
    ],
    outside_value=0.01)
```

DQN

1. replay_buffer

每个时刻存取(obs, act, rew, next_obs, done), 有最大容量, 循环存取. frame存的是uint8,节省内存.

2. critic loss:

```
loss = nn.SmoothL1Loss()
```

Pytorch

1. optimizer中定的是base_lr, learning_rate_scheduler中的学习率是调整的比率, 它和base_lr相乘得到最终的lr.
2. LambdaLR, 传入lambda函数来调整lr

```
# Assuming optimizer has two groups.
lambda1 = lambda epoch: epoch // 30
lambda2 = lambda epoch: 0.95 ** epoch
scheduler = LambdaLR(optimizer, lr_lambda=[lambda1, lambda2])
for epoch in range(100):
    train(...)
    validate(...)
    scheduler.step()
```

3. 增加接口的通用性, 可以传入一些函数类型的参数(lambda函数, 或者回调函数)
4. 把一些常用的函数, 常量统一放入到一个module文件中, 然后用这个module来访问这些函数, 常量. 如pytorch_utils.py中:

```
Activation = Union[str, nn.Module]

_str_to_activation = {
    'relu': nn.ReLU(),
    'tanh': nn.Tanh(),
    'leaky_relu': nn.LeakyReLU(),
    'sigmoid': nn.Sigmoid(),
    'selu': nn.SELU(),
    'softplus': nn.Softplus(),
    'identity': nn.Identity(),
}

def build_mlp(
    input_size: int,
    output_size: int,
    n_layers: int,
    size: int,
    activation: Activation = 'tanh',
    output_activation: Activation = 'identity',
):
```

```
device = None

def init_gpu(use_gpu=True, gpu_id=0):
    global device
    if torch.cuda.is_available() and use_gpu:
        device = torch.device("cuda:" + str(gpu_id))
        print("Using GPU id {}".format(gpu_id))
    else:
        device = torch.device("cpu")
        print("GPU not detected. Defaulting to CPU.")
```

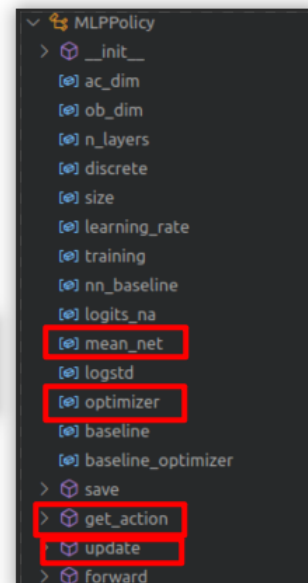
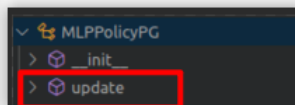
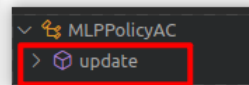
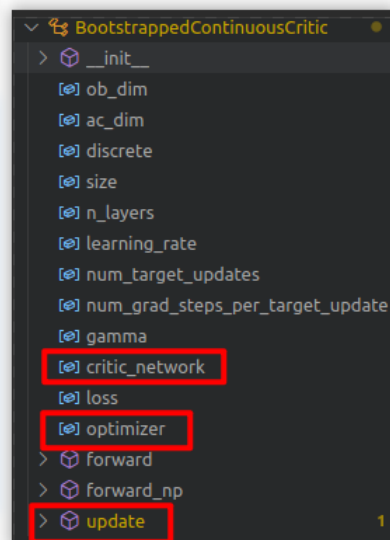
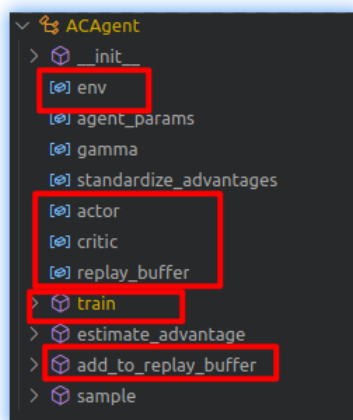
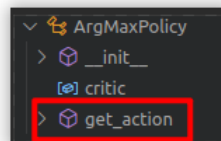
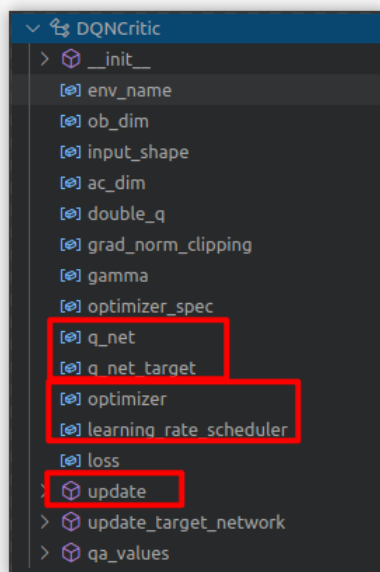
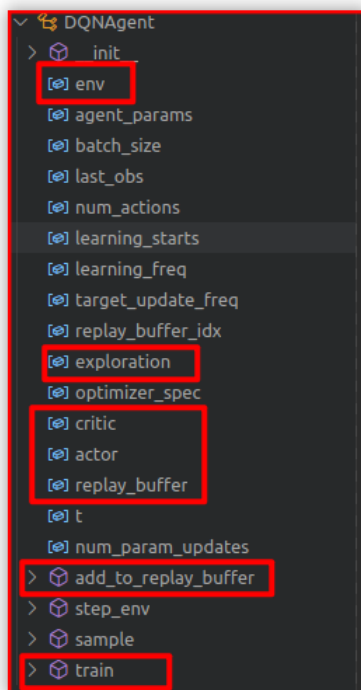
这些方法和变量在不同的项目中都能用到, 而且不用改变.

5. 运行时才sampling

```
idxes = sample_n_unique(lambda: random.randint(0, 10), batch_size)
```

```
def sample_n_unique(sampling_f, n):  
    """Helper function. Given a function `sampling_f` that returns  
    comparable objects, sample n such unique objects.  
    """  
    res = []  
    while len(res) < n:  
        candidate = sampling_f()  
        if candidate not in res:  
            res.append(candidate)  
    return res
```

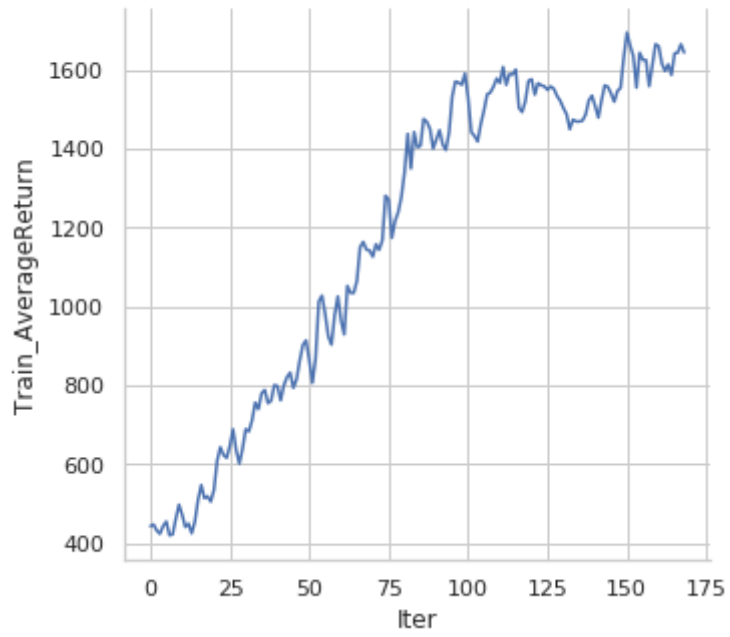
RL framework design



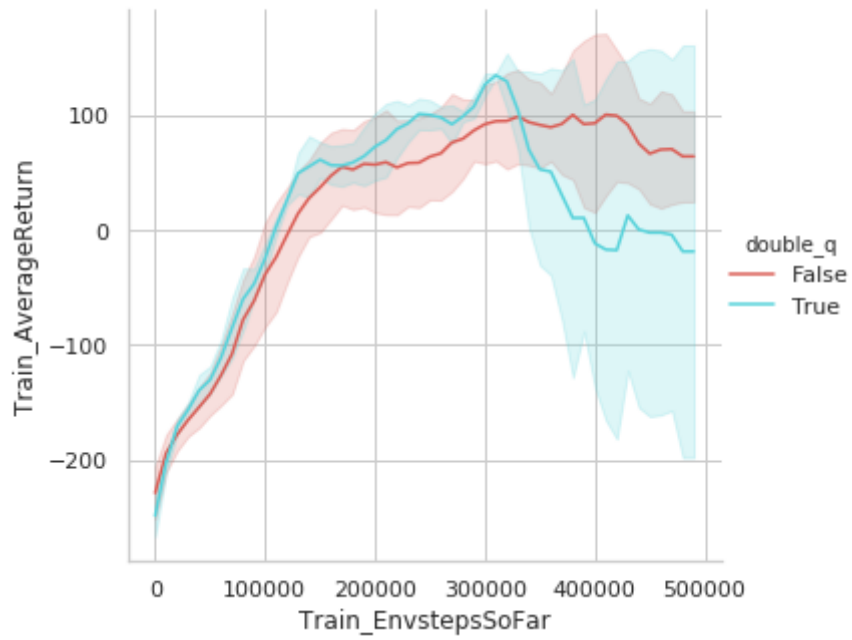
Exps

DQN

q1_dqn_MsPacman-v0



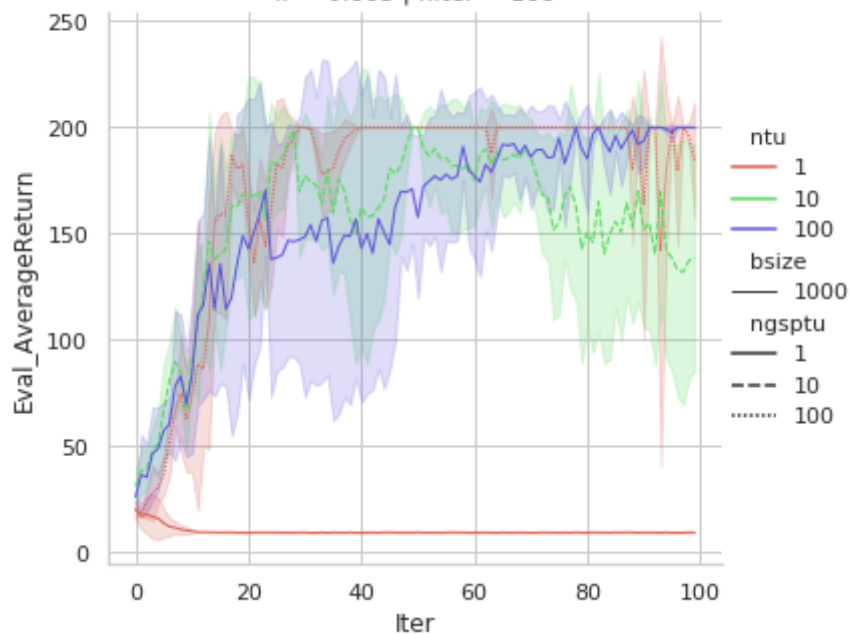
q2_dqn_LunarLander-v3



AC

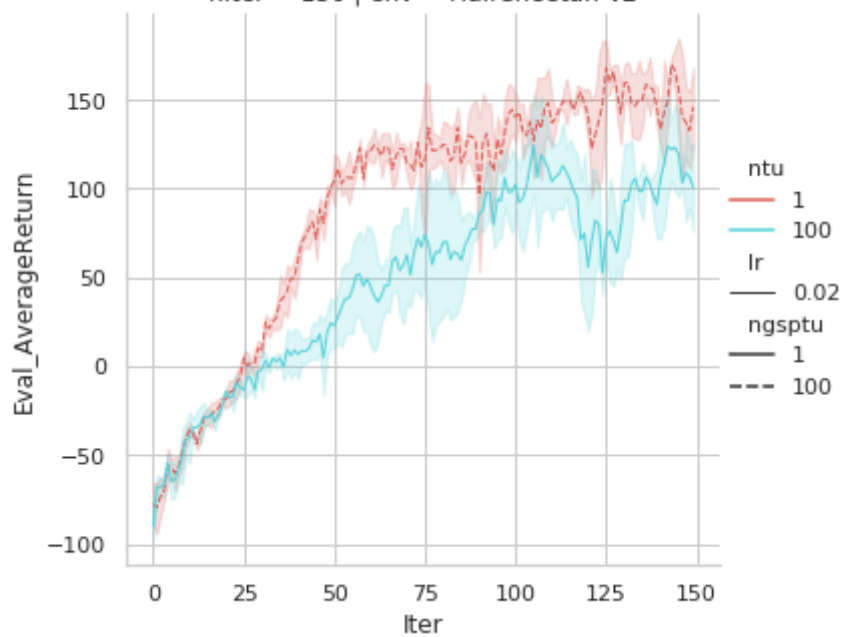
q4_actor-critic_CartPole-v0

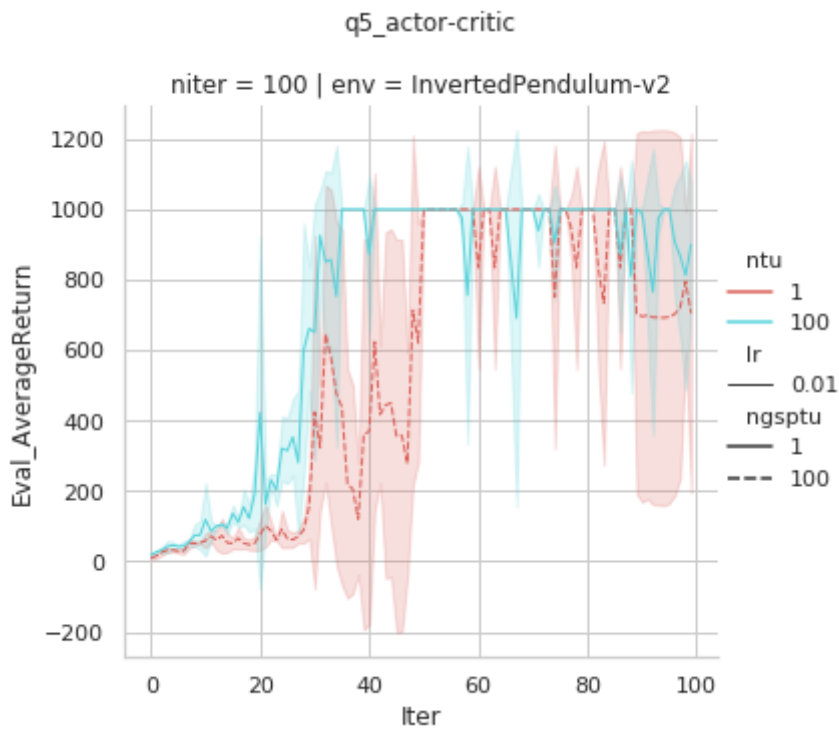
lr = 0.005 | niter = 100



q5_actor-critic

niter = 150 | env = HalfCheetah-v2





```
for _ in range(self.num_target_updates):
    t_n = reward_n + self.gamma * self.forward_np(next_ob_no) * (1.0 - terminal_n)
    for _ in range(self.num_grad_steps_per_target_update):
        ob_no_tensor = ptu.from_numpy(ob_no)
        t_n_tensor = ptu.from_numpy(t_n)
        loss = self.loss(self.forward(ob_no_tensor), t_n_tensor)
        self.optimizer.zero_grad()
        loss.backward()
        self.optimizer.step()
```

ntu: num_target_updates

ngsptu: num_grad_steps_per_target_update

可以看出

1. $ntu \sim ngsptu$, 模型训练效果很差.
2. $ntu \gg ngsptu$, 模型需要花更多时间收敛, 最终效果和 $ngsptu \gg ntu$ 相当
3. $ngsptu \gg ntu$, 模型收敛快, 但是收敛后容易震荡(因为保持target不变的情况下, 进行了更多次的梯度下降)