

# 强化学习基础篇（九）OpenAI Gym基础介绍

## 1. Gym介绍

Gym是一个研究和开发强化学习相关算法的仿真平台，无需智能体先验知识，由以下两部分组成

- Gym开源库：测试问题的集合。当你测试强化学习的时候，测试问题就是环境，比如机器人玩游戏，环境的集合就是游戏的画面。这些环境有一个公共的接口，允许用户设计通用的算法。
- OpenAI Gym服务：提供一个站点和API（比如经典控制问题：CartPole-v0），允许用户对他们的测试结果进行比较。

## 2. Gym安装

我们需要在Python 3.5+的环境中简单得使用pip安装gym

```
1 pip install gym
```

如果需要从源码安装gym，那么可以：

```
1 git clone https://github.com/openai/gym
2 cd gym
3 pip install -e .
```

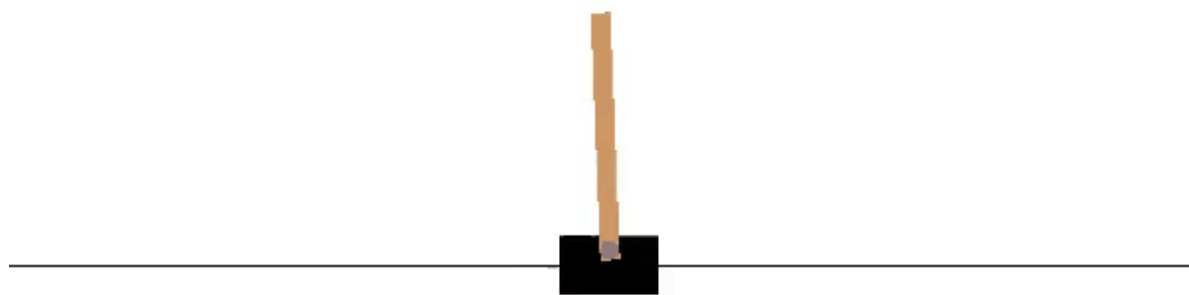
可以运行pip install -e .[all]执行包含所有环境的完整安装。 这需要安装一些依赖包，包括cmake和最新的pip版本。

## 3. Gym使用demo

简单来说OpenAI Gym提供了许多问题和环境（或游戏）的接口，而用户无需过多了解游戏的内部实现，通过简单地调用就可以用来测试和仿真。接下来以经典控制问题CartPole-v0为例，简单了解一下Gym的特点

```
1 # 导入gym环境
2 import gym
3 # 声明所使用的环境
4 env = gym.make('CartPole-v0')
5 # 环境初始化
6 env.reset()
7
8 # 对环境进行迭代执行1000次
9 for _ in range(1000):
10     env.render()
11     observation, reward, done, info = env.step(env.action_space.sample()) #
    采取随机动作
12     if done:
13         env.reset()
14 env.close()
```

运行效果如下



以上代码中可以看出，gym的核心接口是Env。作为统一的环境接口，Env包含下面几个核心方法：

- `reset(self)`: 重置环境的状态，返回观察。
- `step(self, action)`: 推进一个时间步长，返回`observation`, `reward`, `done`, `info`。
- `render(self, mode='human', close=False)`: 重绘环境的一帧。默认模式一般比较友好，如弹出一个窗口。
- `close(self)`: 关闭环境，并清除内存

以上代码首先导入gym库，然后创建CartPole-v0环境，并重置环境状态。在for循环中进行1000个时间步长控制，`env.render()`刷新每个时间步长环境画面，对当前环境状态采取一个随机动作（0或1），在环境返回`done`为True时，重置环境，最后循环结束后关闭仿真环境。

## 4、观测 (Observations)

在上面代码中使用了`env.step()`函数来对每一步进行仿真，在Gym中，`env.step()`会返回 4 个参数：

- 观测 *Observation (Object)*: 当前`step`执行后，环境的观测(类型为对象)。例如，从相机获取的像素点，机器人各个关节的角度或棋盘游戏当前的状态等；
- 奖励 *Reward (Float)*: 执行上一步动作(`action`)后，智能体(*agent*)获得的奖励(浮点类型)，不同的环境中奖励值变化范围也不相同，但是强化学习的目标就是使得总奖励值最大；
- 完成 *Done (Boolean)*: 表示是否需要将环境重置 `env.reset`。大多数情况下，当 *Done* 为True 时，就表明当前回合(*episode*)或者试验(*tial*)结束。例如当机器人摔倒或者掉出台面，就应当终止当前回合进行重置(`reset`)；
- 信息 *Info (Dict)*: 针对调试过程的诊断信息。在标准的智体仿真评估当中不会使用到这个`info`，具体用到的时候再说。

在 Gym 仿真中，每一次回合开始，需要先执行 `reset()` 函数，返回初始观测信息，然后根据标志位 *done* 的状态，来决定是否进行下一次回合。所以更恰当的方法是遵守*done*的标志。

```

1 import gym
2 env = gym.make('CartPole-v0')
3 for i_episode in range(20):
4     observation = env.reset()
5     for t in range(100):
6         env.render()
7         print(observation)
8         action = env.action_space.sample()
9         observation, reward, done, info = env.step(action)
10        if done:
11            print("Episode finished after {} timesteps".format(t+1))
12            break
13    env.close()

```

代码运行结果的片段如下所示：

```

1 [ 0.04025062 -0.04312649  0.00186348  0.02288173]
2 [ 0.03938809 -0.23827512  0.00232111  0.31615203]
3 [ 0.03462259 -0.43343005  0.00864416  0.60956605]
4 [ 0.02595398 -0.23843     0.02083548  0.31961824]
5 [ 0.02118538 -0.43384239  0.02722784  0.6187984 ]
6 [ 0.01250854 -0.23911113  0.03960381  0.33481376]
7 [ 0.00772631 -0.43477369  0.04630008  0.63971794]
8 [-0.00096916 -0.63050954  0.05909444  0.94661444]
9 [-0.01357935 -0.43623107  0.07802673  0.67306909]
10 [-0.02230397 -0.24227538  0.09148811  0.40593731]
11 [-0.02714948 -0.43856752  0.09960686  0.72600415]
12 [-0.03592083 -0.24495361  0.11412694  0.46625881]
13 [-0.0408199  -0.05161354  0.12345212  0.21161588]
14 [-0.04185217  0.14154693  0.12768444 -0.03971694]
15 [-0.03902123 -0.05515279  0.1268901  0.29036807]
16 [-0.04012429 -0.25183418  0.13269746  0.6202239 ]
17 [-0.04516097 -0.05879065  0.14510194  0.37210296]
18 [-0.04633679  0.13400401  0.152544    0.12846047]
19 [-0.04365671 -0.06293669  0.15511321  0.46511532]
20 [-0.04491544 -0.25987115  0.16441551  0.80239106]
21 [-0.05011286 -0.45681992  0.18046333  1.14195086]
22 [-0.05924926 -0.65378152  0.20330235  1.48536419]
23 Episode finished after 22 timesteps

```

上面的结果可以看到这个迭代中，输出的观测为一个列表。这是CartPole环境特有的状态，其规则是  $[x, \theta, \dot{x}, \dot{\theta}]$ 。

其中：

- $x$ 表示小车在轨道上的位置 (*position of the cart on the track*)
- $\theta$ 表示杆子与竖直方向的夹角 (*angle of the pole with the vertical*)
- $\dot{x}$ 表示小车速度 (*cart velocity*)
- $\dot{\theta}$ 表示角度变化率 (*rate of change of the angle*)

## 5、空间 (*Spaces*)

每次执行的动作(action)都是从环境动作空间中随机进行选取的，但是这些动作 (action) 是什么?在 Gym 的仿真环境中，有运动空间 *action\_space* 和观测空间*observation\_space* 两个指标，程序中被定义为 *Space*类型，用于描述有效的运动和观测的格式和范围。下面是一个代码示例：

```

1 import gym
2 env = gym.make('CartPole-v0')
3 print(env.action_space)
4 print(env.observation_space)

```

```

1 Discrete(2)
2 Box(-3.4028234663852886e+38, 3.4028234663852886e+38, (4,), float32)

```

从程序运行结果可以看出：

- `action_space` 是一个离散`Discrete`类型，从`discrete.py`源码可知，范围是一个 $\{0, 1, \dots, n-1\}$  长度为  $n$  的非负整数集合，在`CartPole-v0`例子中，动作空间表示为 $\{0, 1\}$ 。
- `observation_space` 是一个`Box`类型，从`box.py`源码可知，表示一个  $n$  维的盒子，所以在上一节打印出来的`observation`是一个长度为 4 的数组。数组中的每个元素都具有上下界。

## 6. 奖励(reward)

在`gym`的`Cart Pole`环境 (`env`) 里面，左移或者右移小车的`action`之后，`env`会返回一个+1的`reward`。其中`CartPole-v0`中到达200个`reward`之后，游戏也会结束，而`CartPole-v1`中则为500。最大奖励 (`reward`) 阈值可通过前面介绍的注册表进行修改。

## 7. 注册表

`Gym`是一个包含各种各样强化学习仿真环境的大集合，并且封装成通用的接口暴露给用户，查看所有环境的代码如下

```

1 from gym import envs
2 print(envs.registry.all())

```

## 8.注册模拟器

`Gym`支持将用户制作的环境写入到注册表中，需要执行 `gym.make()`和在启动时注册`register`。如果要注册自己的环境，那么假设你在以下结构中定义了自己的环境：

```

1 myenv/
2   __init__.py
3   myenv.py

```

i. `myenv.py` 包含适用于我们自己的环境的类。在 `init.py`中，输入以下代码：

```

1 from gym.envs.registration import register
2 register(
3     id='MyEnv-v0',
4     entry_point='myenv.myenv:MyEnv', # 第一个myenv是文件夹名字，第二个myenv是文件名，MyEnv是文件内类的名字
5 )

```

ii. 要使用我们自己的环境：

```

1 import gym
2 import myenv # 一定记得导入自己的环境，这是很容易忽略的一点
3 env = gym.make('MyEnv-v0')

```

iii. 在PYTHONPATH中安装 myenv 目录或从父目录启动python。

```
1  目录结构:
2  myenv/
3      __init__.py
4      my_hotter_colder.py
5  -----
6  __init__.py 文件:
7  -----
8  from gym.envs.registration import register
9  register(
10     id='MyHotterColder-v0',
11     entry_point='myenv.my_hotter_colder:MyHotterColder',
12 )
13 -----
14 my_hotter_colder.py文件:
15 -----
16 import gym
17 from gym import spaces
18 from gym.utils import seeding
19 import numpy as np
20
21 class MyHotterColder(gym.Env):
22     """Hotter Colder
23     The goal of hotter colder is to guess closer to a randomly selected
24     number
25
26     After each step the agent receives an observation of:
27     0 - No guess yet submitted (only after reset)
28     1 - Guess is lower than the target
29     2 - Guess is equal to the target
30     3 - Guess is higher than the target
31
32     The rewards is calculated as:
33     (min(action, self.number) + self.range) / (max(action, self.number) +
34     self.range)
35
36     Ideally an agent will be able to recognise the 'scent' of a higher
37     reward and
38     increase the rate in which is guesses in that direction until the reward
39     reaches
40     its maximum
41     """
42     def __init__(self):
43         self.range = 1000 # +/- value the randomly select number can be
44         between
45         self.bounds = 2000 # Action space bounds
46
47         self.action_space = spaces.Box(low=np.array([-self.bounds]),
48         high=np.array([self.bounds]))
49         self.observation_space = spaces.Discrete(4)
50
51         self.number = 0
52         self.guess_count = 0
53         self.guess_max = 200
54         self.observation = 0
```

```

50         self.seed()
51         self.reset()
52
53     def seed(self, seed=None):
54         self.np_random, seed = seeding.np_random(seed)
55         return [seed]
56
57     def step(self, action):
58         assert self.action_space.contains(action)
59
60         if action < self.number:
61             self.observation = 1
62
63         elif action == self.number:
64             self.observation = 2
65
66         elif action > self.number:
67             self.observation = 3
68
69         reward = ((min(action, self.number) + self.bounds) / (max(action,
self.number) + self.bounds)) ** 2
70
71         self.guess_count += 1
72         done = self.guess_count >= self.guess_max
73
74         return self.observation, reward[0], done, {"number": self.number,
"guesses": self.guess_count}
75
76     def reset(self):
77         self.number = self.np_random.uniform(-self.range, self.range)
78         self.guess_count = 0
79         self.observation = 0
80         return self.observation

```

## 9. OpenAI Gym评估平台

用户可以记录和上传算法在环境中的表现或者上传自己模型的Gist，生成评估报告，还能录制模型玩游戏的小视频。在每个环境下都有一个排行榜，用来比较大家的模型表现。

上传于录制方法如下所示

```

1  import gym
2  from gym import wrappers
3  env = gym.make('CartPole-v0')
4  env = wrappers.Monitor(env, '/tmp/cartpole-experiment-1')
5  for i_episode in range(20):
6      observation = env.reset()
7      for t in range(100):
8          env.render()
9          print(observation)
10         action = env.action_space.sample()
11         observation, reward, done, info = env.step(action)
12         if done:
13             print("Episode finished after {} timesteps".format(t+1))
14             break

```

使用Monitor Wrapper包装自己的环境，在自己定义的路径下将记录自己模型的性能。支持将一个环境下的不同模型性能写在同一个路径下。

在[官网](#)注册账号后，可以在个人页面上看到自己的API\_Key，接下来可以将结果上传至OpenAI Gym：

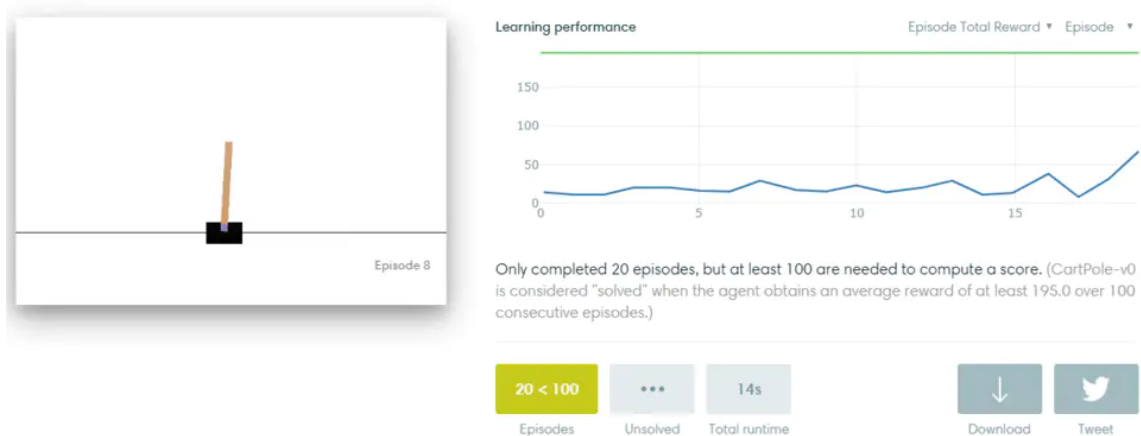
```
1 import gym
2 gym.upload('/tmp/cartpole-experiment-1', api_key='YOUR_API_KEY')
```

然后得到如下结果：

```
*****
You successfully uploaded your evaluation on CartPole-v0 to
OpenAI Gym! You can find it at:

https://gym.openai.com/evaluations/eval\_hxhpHLIqRaGT2xyN3KaAoQ
*****
```

打开链接会有当前模型在环境下的评估报告，并且还录制了小视频：



每次上传结果，OpenAI Gym都会对其进行评估。

Nice work on your evaluation! For your result to be reviewed, please [provide a Gist](#) explaining how to reproduce your results. Please read our [evaluation guidelines](#) for some background, or [pop into chat](#) with any questions.

Gist URL

[Attach Gist](#)

创建一个Github Gist将结果上传，或者直接在upload时传入参数：

```
1 import gym
2 gym.upload('/tmp/cartpole-experiment-1',
  writeup='https://gist.github.com/gdb/b6365e79be6052e7531e7ba6ea8caf23',
  api_key='YOUR_API_KEY')
```

评估将自动计算得分，并生成一个漂亮的页面。

在大多数环境中，我们的目标是尽量减少达到阈值级别的性能所需的步骤数。不同的环境都有不同的阈值，在某些环境下，尚不清楚该阈值是什么，此时目标是使最终的表现最大化。在cartpole这个环境中，阈值就是立杆能够直立的帧数。