# 1 项目文件结构

- `root`：
  - `README`
    - `README.md`
  - `log_plotter.py`：用于：
    - 将 *value-based* 方法中的 *Dueling DDQN* 和 *Dueling DQN* 模型在训练过程中的 `log_score.txt` 和 `log_loss.txt` 绘制成变化曲线，存储在 `./output/value_based/` 目录下；
    - 将 *policy-based* 方法中的 *TD3* 和 *DDPG* 模型在 *Humanoid-v2* 环境下训练过程中的 *avg reward in evaluation* 绘制成变化曲线，存储在 `./output/policy_based/` 目录下；
  - `output`：
    - `value_based`：存储 *log_plotter.py* 绘制的 *value-based* 相关的图像
    - `policy_based`：存储 *log_plotter.py* 绘制的 *policy-based* 相关的图像
  - `value_based`：
    - `DuelingDQN` / `DuelingDDQN`：
      - `log`：用于记录训练过程中的 *score* 和 *loss* 等信息
      - `model`：用于保存训练得到的 `Dueling_DQN_breakout.pkl` / `Dueling_DDQN_breakout.pkl` 模型
      - `atari_playground.py`：主函数
      - `RL_brain.py`：实现 `Dueling_DeepQNetwork` / `Dueling_Double_DeepQNetwork` 类
      - `atari_wrappers.py`：*Atari* 环境的封装函数
      - `utils.py`：用于模型和训练信息保存的工具函数
  - `policy_based`：
    - `TD3_results`：保存了四种 *MuJoCo* 环境下的训练信息、绘制的曲线以及模型
      - `human`：
        - `TD3` / `DDPG`：保存 *TD3* / *DDPG* 算法的模型结果
          - `models`：模型
          - `output`：输出的曲线图片
          - `results`：输出的用于绘制曲线的数组
      - `ant` / `halfcheetah` \ `hopper`：
        - `models`：模型
        - `output`：输出的曲线图片
        - `results`：输出的用于绘制曲线的数组
    - `TD3`：*policy-based* 模型的主文件夹
      - `main.py`：主函数
      - `TD3.py`：*TD3* 类的实现
      - `DDPG.py`：*DDPG* 类的实现
      - `utils.py`：*ReplayBuffer* 类的实现
      - `models`：模型
      - `output`：输出的曲线图片

- **results**：输出的用于绘制曲线的数组

---

# 2 项目使用方法

## 2.1 *value-based*

- 假设现在想要运行 *Dueling DDQN* 模型，*Dueling DQN* 模型同理
- 首先进入所创建的 `conda python3.7` 虚拟环境中，进入到 `./value_based/Dueling_DDQN/` 目录下，使用如下指令开始：

```
1  python atari_playground.py --env_name BreakoutNoFrameskip-v4
```

- 运行过程中的 *evalueate* 会周期性输出三个值，分别是：*evaluate episode*、该 *evaluate* 周期内的平均 *reward*，以及该 *evaluate* 周期内的最大 *reward*
- 模型的状态 `dict` 被保存，因此若想测试现有模型，可以直接将 `atari_playground.py` 中的模型初始化替换为模型的加载

## 2.2 *policy-based*

- 首先进入所创建的 `conda python3.7` 虚拟环境中，进入到 `./policy_based/TD3/` 目录下，使用如下指令开始，该方法会使用 *TD3* 算法在 *Humanoid-v2* 上进行训练，并保存模型；

```
1  python main.py --env Humanoid-v2 --policy TD3 --save_model
```

- 若想要验证已有模型的效果，请添加 `--load_model default` 字段，该模型会自动去默认文件夹下寻找并加载已有模型。具体的命令行参数设置如下所示：

```
1  parser = argparse.ArgumentParser()
2      # Policy name (TD3 or DDPG)
3      parser.add_argument("--policy", default="TD3")
4      # OpenAI gym environment name
5      parser.add_argument("--env", default="HalfCheetah-v2")
6      # Sets Gym, PyTorch and Numpy seeds
7      parser.add_argument("--seed", default=0, type=int)
8      # Time steps initial random policy is used
9      parser.add_argument("--start_timesteps", default=2e3, type=int)
10     # How often (time steps) we evaluate
11     parser.add_argument("--eval_freq", default=5e3, type=int)
12     # Max time steps to run environment
13     parser.add_argument("--max_timesteps", default=1e6, type=int)
14     # Std of Gaussian exploration noise
15     parser.add_argument("--expl_noise", default=0.1)
16     # Batch size for both actor and critic
17     parser.add_argument("--batch_size", default=256, type=int)
```

```
18        # Discount factor
19        parser.add_argument("--discount", default=0.99)
20        # Target network update rate
21        parser.add_argument("--tau", default=0.005)
22        # Noise added to target policy during critic update
23        parser.add_argument("--policy_noise", default=0.2)
24        # Range to clip target policy noise
25        parser.add_argument("--noise_clip", default=0.5)
26        # Frequency of delayed policy updates
27        parser.add_argument("--policy_freq", default=2, type=int)
28        # Save model and optimizer parameters
29        parser.add_argument("--save_model", action="store_true")
30        # Model load file name, "" doesn't load, "default" uses file_name
31        parser.add_argument("--load_model", default="")
32        args = parser.parse_args()
```

注意到，当我们选择加载已有模型后，整个模型能够在 `learn_start` 之后立刻达到很高的 *reward*，如下图所示：