# 工作总结报告

## 目标

使用DDPG算法对人形 双足机器人模型进行训练，使其学会直立行走

## 实现过程

在openAI上baselines包中有完整的强化学习的算法包，如：DQN、PPO、DDPG，通过设定相关的参数，选择使用 DDPG算法，配置Humanoid-v3环境(Mujoco)，开始训楼，训练结束后，展示当前数据量训练所得模型的结果。

- 具体使用方法：

```
xinSamdeMacBook-Pro:baselines-master sam$ python3 -m baselines.run -h
usage: run.py [-h] [--env ENV] [--env_type ENV_TYPE] [--seed SEED] [--alg ALG]
              [--num_timesteps NUM_TIMESTEPS] [--network NETWORK]
              [--gamestate GAMESTATE] [--num_env NUM_ENV]
              [--reward_scale REWARD_SCALE] [--save_path SAVE_PATH]
              [--save_video_interval SAVE_VIDEO_INTERVAL]
              [--save_video_length SAVE_VIDEO_LENGTH] [--log_path LOG_PATH]
              [--play]

optional arguments:
  -h, --help            show this help message and exit
  --env ENV             environment ID (default: Reacher-v2)
  --env_type ENV_TYPE   type of environment, used when the environment type
                        cannot be automatically determined (default: None)
  --seed SEED           RNG seed (default: None)
  --alg ALG             Algorithm (default: ppo2)
  --num_timesteps NUM_TIMESTEPS
  --network NETWORK     network type (mlp, cnn, lstm, cnn_lstm, conv_only)
                        (default: None)
  --gamestate GAMESTATE
                        game state to load (so far only used in retro games)
                        (default: None)
  --num_env NUM_ENV     Number of environment copies being run in parallel.
                        When not specified, set to number of cpus for Atari,
                        and to 1 for Mujoco (default: None)
  --reward_scale REWARD_SCALE
                        Reward scale factor. Default: 1.0 (default: 1.0)
  --save_path SAVE_PATH
                        Path to save trained model to (default: None)
  --save_video_interval SAVE_VIDEO_INTERVAL
                        Save video every x steps (0 = disabled) (default: 0)
  --save_video_length SAVE_VIDEO_LENGTH
                        Length of recorded video. Default: 200 (default: 200)
  --log_path LOG_PATH   Directory to save learning curve data. (default: None)
  --play
```

- Episodes = 1000 时的 reward

- Episodes = 10000 时的 reward
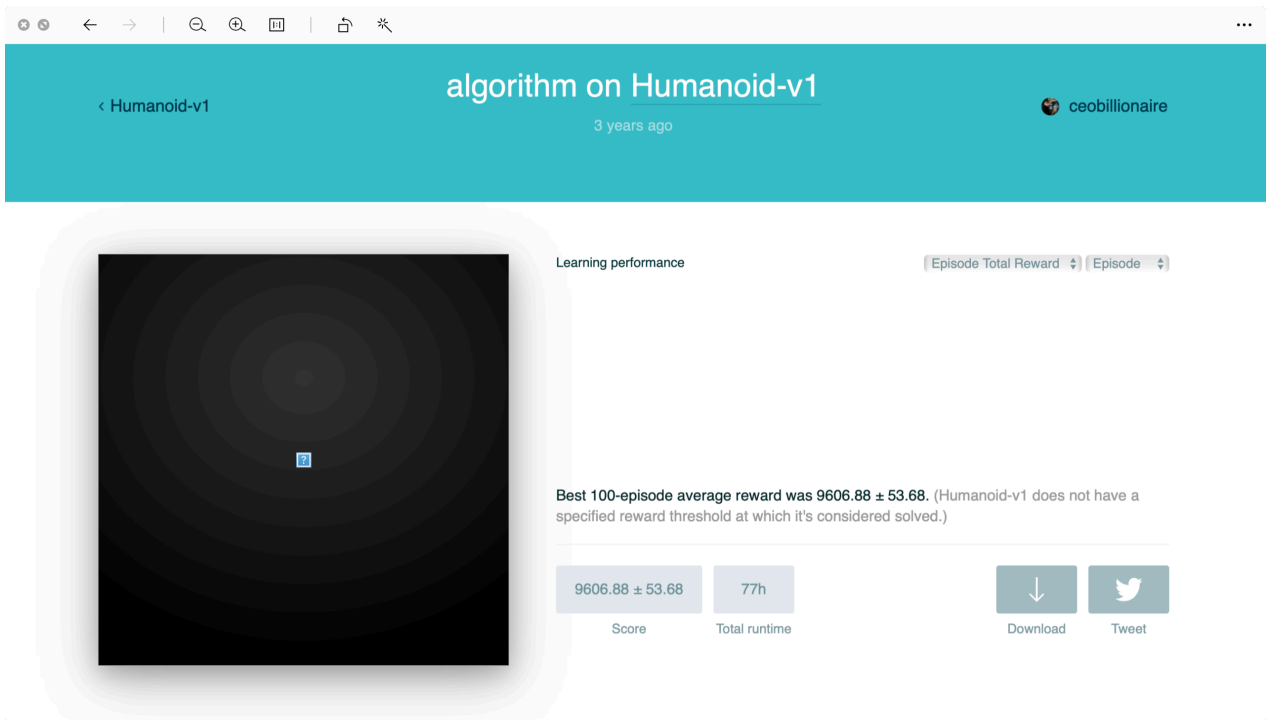


- Episodes = 100000 时的 reward

可看出，随着训练的数据量不断增加，所得到的reward的均值也在提升，达到了训练的目的。经查询资料得知，需要训练77h，平均reward达到9000分上下，才可完成对走路的学习。

# algorithm on Humanoid-v1

3 years ago

ceobillionaire

Learning performance

Episode Total Reward ⇅    Episode ⇅

**Best 100-episode average reward was 9606.88 ± 53.68.** (Humanoid-v1 does not have a specified reward threshold at which it's considered solved.)

| 9606.88 ± 53.68 | 77h |
|---|---|
| Score | Total runtime |

↓
Download

Tweet

# 心得体会

经过这次实习，我们深入了解了协作完成一个项目的方法与技巧。并且对新兴的人工智能领域有了一定的认识。同时，我们还调整了自己在学校时混乱的作息，坚持早睡早起，切身体会到了通勤的感觉。可以说是收获颇丰，希望今后的学习生活中能够运用到本次实习所学到的知识和技巧，继续成长。