

PPO × Family 第二讲技术问题 QA

课程组整合梳理了自第二节课发布以来的相关技术问题答疑，希望能启发更多对于决策智能和强化学习的思考，也欢迎大家贡献问题或者积极参与讨论~

Q0: 有没有第二节课内容的大白话总结？

A0: 决策问题环境类型的多样性带来了不同类型的动作空间，而处理不同的动作空间需要不同的策略建模方式：

小节	技术要点	代表决策任务
离散动作空间	<ul style="list-style-type: none">离散动作空间计算图 + logit构建概率分布及采样动作	火箭回收 视频游戏（Atari, Procgen）
连续动作空间	<ul style="list-style-type: none">连续动作空间计算图/梯度流PPO 和 DDPG 在设计理念和数据属性两方面的对比	机器人控制（MuJoCo） 无人机姿态控制
多维离散空间	<ul style="list-style-type: none">multi-discrete 等四种网络结构的对比multi-discrete 和 discrete 在探索和利用角度的分析	交通信控 键盘操作的游戏
混合动作空间	<ul style="list-style-type: none">各类动作空间离散化方法Hybrid PPO (HPPO)HPPO + 自回归/分离预测器	导航控制 即时战略游戏（星际争霸2） MOBA 类游戏（DOTA2）

总之，对一个实际的决策问题环境选用适合的建模方式，转化为标准的 MDP 形式，是解决决策问题最关键的步骤。所谓，阵而后战，兵法之常；运用之妙，存乎一心。

Q1: 第二节课代码题的第一题，就是在给出的代码段里（即 `reparam_grad` 函数），根据重参数化方法对应得到的梯度公式，用 `numpy` 实现公式，并运行整个程序进行对比，这道题是这么理解的吗？

A1: 这道题就是对应梯度公式实现相应逻辑，可以参考代码段上面的 `naive_grad`，这个函数就是实现了朴素方法的梯度公式。整体这道题的目的是希望通过这个例子来对比解释重参数化到底是如何减小梯度方差，通过直观的数据对比来说明这一点。

Q2: 第二节课的应用 demo 都挺有意思的，但是实际运行起来配环境时出现了不少bug，并且究竟运行成什么样才算效果好，作为初学者不太有概念，这方面课程官方能提供一些帮助吗？

A2: 课程组也注意到了这个问题，对于第二节课的题解和第三节课的实践题，我们会做以下改进：

- 对于环境安装，课程组会给出完整的安装指南 + 可以直接运行环境 demo 的 docker + 可以直接运行环境 demo 的在线 notebook（类似 colab 和 kaggle 上的 notebook）
- 对于训练过程，课程组会从第三节课作业起预先给出完整的训练日志和相应的操作录屏，帮助大家建立一个更直观的认识

- 对于代码和使用到的一些开源库，我们会安排一组系列小视频来讲解相关的代码（包括但不限于 gym 环境定义，DI-engine 代码详解，DI-treetensor 和 wandb 使用指南等等）

Q3: 请问一下，课程2中的辅助材料，也是需要学习的吗？这些材料和动作空间的关系是？

A3: 辅助材料是作为补充使用的，根据个人兴趣和时间安排决定即可。不过课程组建议，就算不深入看也应该简单读一下了解大概内容，以后如果遇到类似的问题可以借助补充材料拓展知识。对于第二节的补充材料，主要对应三部分的内容：

1. 重参数化部分详解连续动作空间输出形式的设计，并对比 PPO 和 SAC 在这方面的异同点。
2. DDPG 对比部分主要深入解释 PPO 中重要性采样（IS）的设计动机，虽然都是做连续动作控制但是两者优化思路不同，因此设计就大相径庭，进而实践中 DDPG 和 PPO 的超参数设置方式也有很大区别。
3. HyAR 部分则是给出了动作空间表征学习的一个例子，这也是目前复杂动作空间研究的主要方向。

Q4: 想问下课上的内容，为什么可以用 Behavior Cloning（BC）来测试网络架构利用数据能力的高低呢？课程里面说的这种专家数据的 BC 具体是怎么做的呢？是用专家数据拟合程度给 reward 么？还是说用 PPO 的网络直接做监督学习呢？

A4: 这里是使用了 PPO 应用于 multi-discrete action space 的两种神经网络模型架构，一种是标准的 discrete 形式，另一种是 multi-discrete。不同模型使用相同的专家数据去做监督学习，实验中发现 multi-discrete 形式的网络学习效果更好，说明这种网络结构设计更有优势。

Q5: 请问 DPPO (Distributed Proximal Policy Optimization) 在更新模型的时候，优势值一般要不要做归一化呢？例如有看到过类似下方的代码实现：

```
1 adv = (adv - adv.mean()) / (adv.std() + 1e-6) # sometimes helpful
```

并且，如果要做归一化，应该分批次做吧，因为 DPPO 每次训练的数据来自于多个 worker，应该是在各自 worker 内产生的数据进行归一化吧？不然就混乱了。

A5: 优势值归一化（adv norm）需要考虑 reward 的数值范围：

- 如果 reward 绝对值在 0-100 之内其实影响不大；
- 如果绝对值大于这个范围，且 reward 波动的确实很明显（注意要和稀疏 reward 区分），那适合用 adv norm。然后这种比较直接的 adv norm，应该 batch 越大统计量越准，所以在实现中，像 DPPO 的话，应该是训练的多卡之间 allreduce 同步 mean 和 std，然后再 norm。

Q6: 加 adv norm 对解决训练过程中出现 NaN 导致崩溃有帮助吗？

A6: 出现 NaN 有很多原因，其中之一就是 adv norm；具体情况是，如果数据多样性很差，一个 batch 里的数据太相近太过相似，那么算出来的 std 就很接近于0，这样 norm 操作时一除就导致 NaN。

Q7: 在 DPPO 中添加 adv norm，有两种实现方式：

- 类似课程组之前将的 allreduce [可参考博客] 出各个 worker 所有的 adv，然后求 mean 和 std，
- 每个 batch 的数据中，是由不同的 worker 贡献的，单独在各自 worker 贡献的数据内求 mean 和 std，然后对各自 worker 贡献的数据进行归一化，

这两个哪个好呢？我自己做了下实验，发现这两种没啥差别。

A7: mean 和 std 计算用的样本肯定是越多越准，但因为 RL 本身数据分布就一直在变，所以可能有些场景里对最终性能影响不大，就跟你的实验结果一样。你要想真正深究这个问题，应该要去可视化这两种设定下算出来的 mean 和 std 的变化情况，再分析这个变化对于智能体性能的影响，并在不同类型的环境上做对比看能不能找到普适结论；决策问题（环境）之间的差异性太大了，所以经验性结论经常变化，但是分析方法和手段掌握了之后，具体问题具体分析就好，没有什么玄学。

Q8: 想问下 PPO 的连续动作的方差，一般是用神经网络训练输出，还是给固定值（或者从固定值开始线性衰减）？我自己试验了一下，发现方差给固定值居然要比用神经网络学出方差要好。对此我 google 了一下，发现 OpenAI 的开源库中使用 PPO 算法时，居然说他们使用了固定的连续动作的方差 sigma，而不是学出来的，链接如下：

https://spinningup.openai.com/en/latest/spinningup/rl_intro.html

A8: PPO 连续动作的 sigma，其实在不同版本的实现里一共有三种，根据不同情况使用：

- fixed：固定 sigma，常用于一些特殊控制任务，如果对环境的 sigma 的设置有足够的先验知识可以这样做
- independent：即为一个可优化的网络参数，但是和 state 无关，是一个独立参数。这是一般 PPO 常用的情形
- state conditioned：由 state 输入通过一定的网络层生成，这种情况在 SAC 中常用，PPO 中较少见。不过有 paper 在 MuJoCo 环境上做过对比实验，至少在 MuJoCo 这样的控制环境上差别不大

三种类型的代码对比可以参考这里的代码: <https://github.com/opendilab/DI-engine/blob/main/ding/model/common/head.py#L965>

中文版的注释详解可以看这个: https://opendilab.github.io/PPOxFamily/continuous_zh.html