

Rapport DDPG

Simon Groc

March 9, 2025

Contents

1	Introduction	1
2	Marche aléatoire	1
3	Calcule de la loss du critique	2
4	Calcule de la loss de l'actor	3

1 Introduction

Le but de ce document est de clarifier le fonctionnement de ddpq au travers d'un parcours du code BBRL que nous avons fait avec des explications du cours et la retranscription de celle-ci en code python.

2 Marche aléatoire

Avant de commencer à entraîner le modèle il nous faut des données sur l'environnement pour cela nous allons devoir nous déplacer de manière aléatoire dans celui-ci et récupérer des échantillons de la forme suivante ("env/terminated", "env/reward", "critic/q_value", "target-critic/q_value"). Ces informations sont récupérées à chaque étape de la marche aléatoire dans l'environnement dans notre cas de figure l'on en fait 10000 avec un agent qui utilise un Gaussian pour se déplacer.

```
class AddGaussianNoise(Agent):
    def __init__(self, sigma):
        super().__init__()
        self.sigma = sigma

    def forward(self, t, **kwargs):
        act = self.get("action", t)
        dist = Normal(act, self.sigma)
        action = dist.sample()
        self.set("action", t, action)
```

3 Calcul de la loss du critique

le calcul de la loss du critique et une fonctions qui prent les parametre suivent:

- `cfg`: Le dictionner des parametre du modelle
- `reward`: Un tableaux des reconpense en fonctions des transition (taille 2xB)
- `must_bootstrap`: Un tableaux de booléen qui indique la fin d'un episode (taille 2xB)
- `q_values`: Le calcul de la Q-values pare le critique (taille 2xB)
- `target_q_values`: La Q-values calculer par le target (taille 2xB)

ensuite on dois calculer cette formule $y_i = r_{i+1} + \gamma \hat{Q}_\phi^{\pi^\theta}(s_{i+1}, \pi(s_{i+1}))$

pour pouvoir la comparer a $\hat{Q}_\phi^{\pi^\theta}(s_i, a_i | \phi)$

et caculer la loss $L = \frac{1}{N} \sum_i \left(y_i - \hat{Q}_\phi^{\pi^\theta}(s_i, a_i | \phi) \right)^2$

Pour se faire on commense par le cacule de $\hat{Q}_\phi^{\pi^\theta}(s_i, a_i | \phi)$. Il nous faut donc indique a l'agent de caculer les nouvelle q_values pour sa il va utilise les état et les actions du workspace trouver pendands la marche aléatoire puis recuper les q_values du workspace.

```
q_agent(rb_workspace, t=0, n_steps=1)
q_values = rb_workspace["critic/q_values"]
```

une fois que on a cette valeur on cacule notre q_valeus courente il nous faux $\hat{Q}_\phi^{\pi^\theta}(s_{i+1}, \pi(s_{i+1}))$ et pour sa on comment pare executer le actore sur s_{i+1} se qui nous donne les action maximale $\pi(s_{i+1})$ une fois optenus on peux enfin metre a jour les valeur du target critique avec s_{i+1} comme indique dans la formule $\hat{Q}_\phi^{\pi^\theta}(s_{i+1}, \pi(s_{i+1}))$.

```
with torch.no_grad():
```

```
ag_actor(rb_workspace, t=1, n_steps=1)
```

```
target_q_agent(rb_workspace, t=1, n_steps=1, detach_actions=True)
```

```
post_q_values = rb_workspace["target-critic/q_values"]
```

Dans le bloc ci-dessu aucun gradient nais pris en conte.

Ensuite on appelle la loss fonction :

```
target = reward[: -1].squeeze() + cfg.algorithm.discount_factor * target_q_values[1].squeeze()
```

```
mse = nn.MSELoss()
```

```
critic_loss = mse(q_values[0].squeeze(-1), target)
```

On commense par calculer $y_i = r_{i+1} + \gamma \hat{Q}_\phi^{\pi^\theta}(s_{i+1}, \pi(s_{i+1}))$ dans cette formule r_{i+1} et `reward[: -1]`, γ represent le degrés d'atténuation du modelle, et $\hat{Q}_\phi^{\pi^\theta}(s_{i+1}, \pi(s_{i+1}))$ c'est simplement q_valeus de l'état suivent si on utilise la meilleur actions pour cette état tous sa est multiplier par `must_bootstrap[1].int()` un tableaux de booléen qui indique si un episode est finis de cette manier on ne calcule pas la fin d'un episode et le debut d'un nouveaux comme une transitions.

`torch.squeeze()` : cette methode supprime toute les dimensions de taille 1.
puis on applique un descente de gradient classique.

4 Calcul de la loss de l'actor

Pour calculer la loss de l'actor il nous faut recalculer les meilleures actions puis à nouveau le critique car on veut récupérer les q_values que on va passer au en paramètre de la loss fonction.

```
ddpg.t_actor(rb_workspace, t=0, n_steps=1)
ddpg.t_critic(rb_workspace, t=0, n_steps=1)

q_values = rb_workspace["critic/q_value"]
```

puis on envoie les q_values dans la loss fonction $\mathbb{E}_{a_t \sim \pi_\theta(\cdot)} [\nabla_a \hat{Q}_\phi^{\pi_\theta}(s_t, a_t) \nabla_\theta \pi_\theta(s_t)]$
pour calculer celle là c'est très simple il suffit de faire moins la moyenne car nous cherchons à la maximiser est pytorch ne fait pas de descente de gradient que pour des problèmes de minimisations.