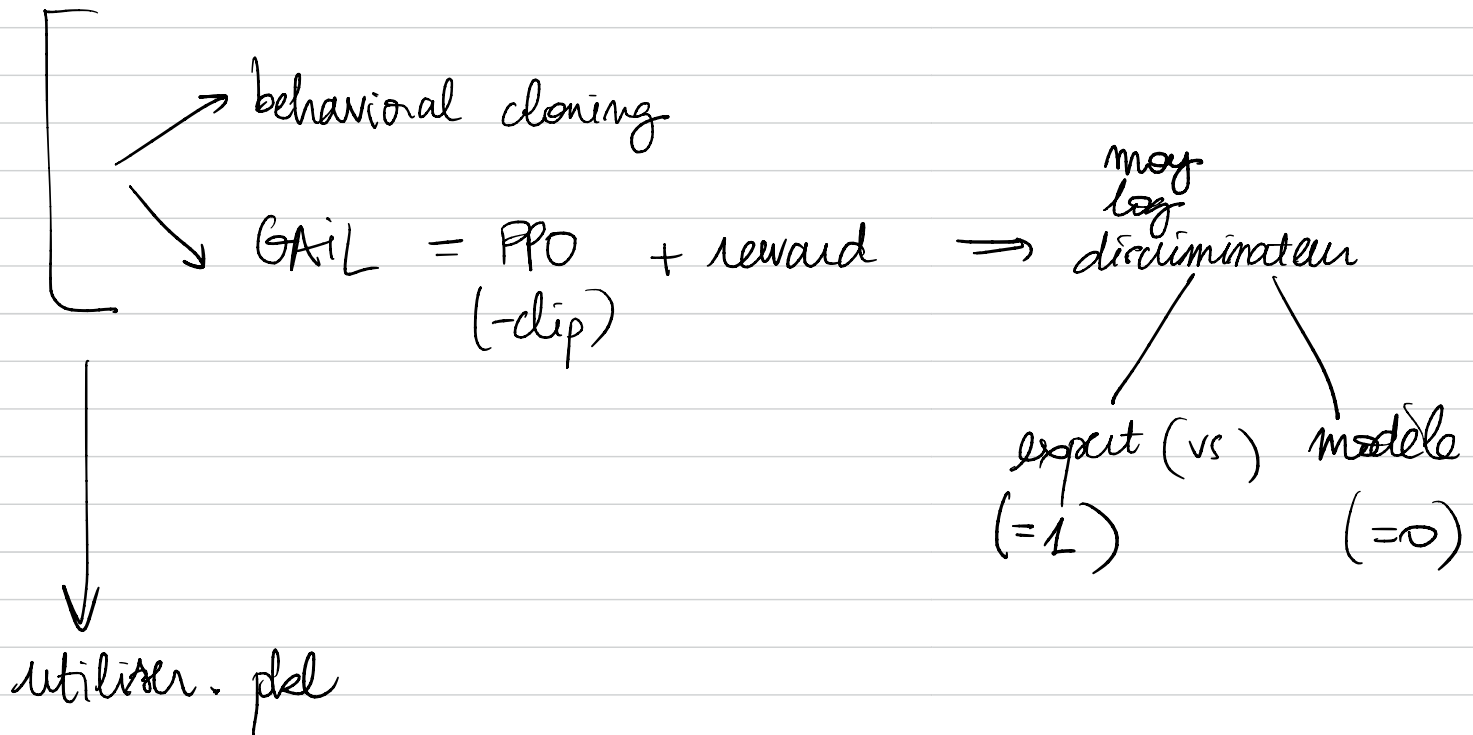
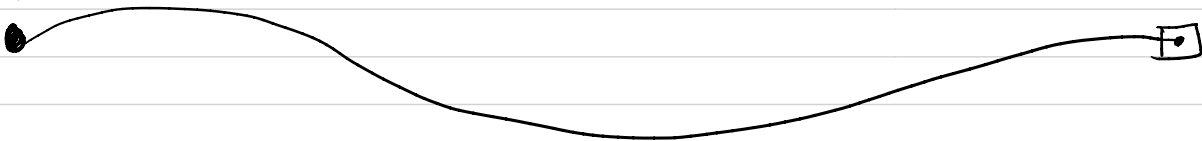


TP12 - Imitation Learning

→ in LunarLander



torch.nn.functional.one-hot.

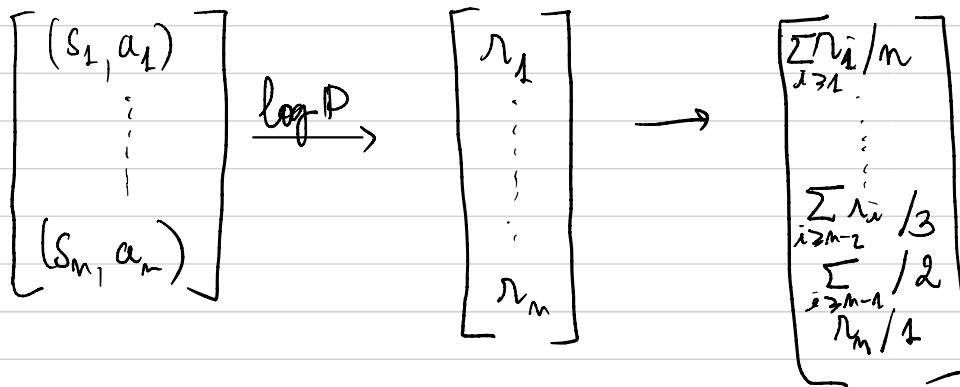


Behavioral learning: apprenticeship supervised sur $\left\{ \begin{matrix} (s_1, a_1) \\ \vdots \\ (s_n, a_n) \end{matrix} \right\}$

But: $\max_{\theta} \sum_{(a_i, s_i) \in \text{expert}} \log \pi_{\theta}(a_i | s_i)$

$\mathbb{E}_{s, a \sim \text{traj}(s_0, a_0)} \left[\log P(\text{expert} | s, a) | s_0, a_0 \right]$

→ Entraînement du discriminateur: étiquetter artificiellement les (s, a) générées / réelles.



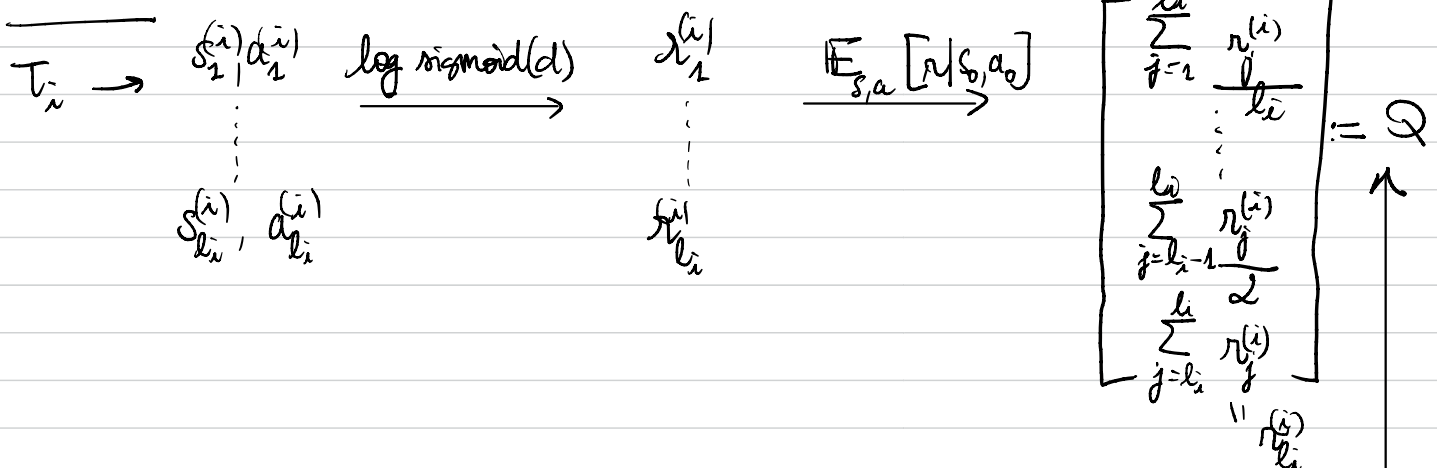
$t_1, \dots, t_m \sim \pi_\theta$ (trajectoires selon π_θ , politique actuelle)
 où $t_i = ((s_1^{(i)}, a_1^{(i)}), \dots, (s_{l_i}^{(i)}, a_{l_i}^{(i)}))$ où l_i longueur de la traj. i .

① Entraîner le discriminateur $d \in \mathbb{R}$ (logit):

$$\mathbb{E}_{(s,a) \sim \tilde{E}_{\text{out}}} [\log(d(s,a))] + \mathbb{E}_{(s,a) \sim \tau} [\log(1-d(s,a))]$$

$$= \text{BCE With logit loss} \left(\underset{E \sim \tilde{E}_{\text{out}}}{(s,a)}, 1 \right) + \text{BCE With logit loss} \left(\underset{\substack{E \sim \tau \\ \text{(pol. actuelle)}}}{(s,a)}, 0 \right)$$

② Rewards:



③ PPO avec la fonction Q .