

Policy Gradient

- David Silver's Slides presented by Jonathan Lorraine
- Images in this powerpoint provided by:

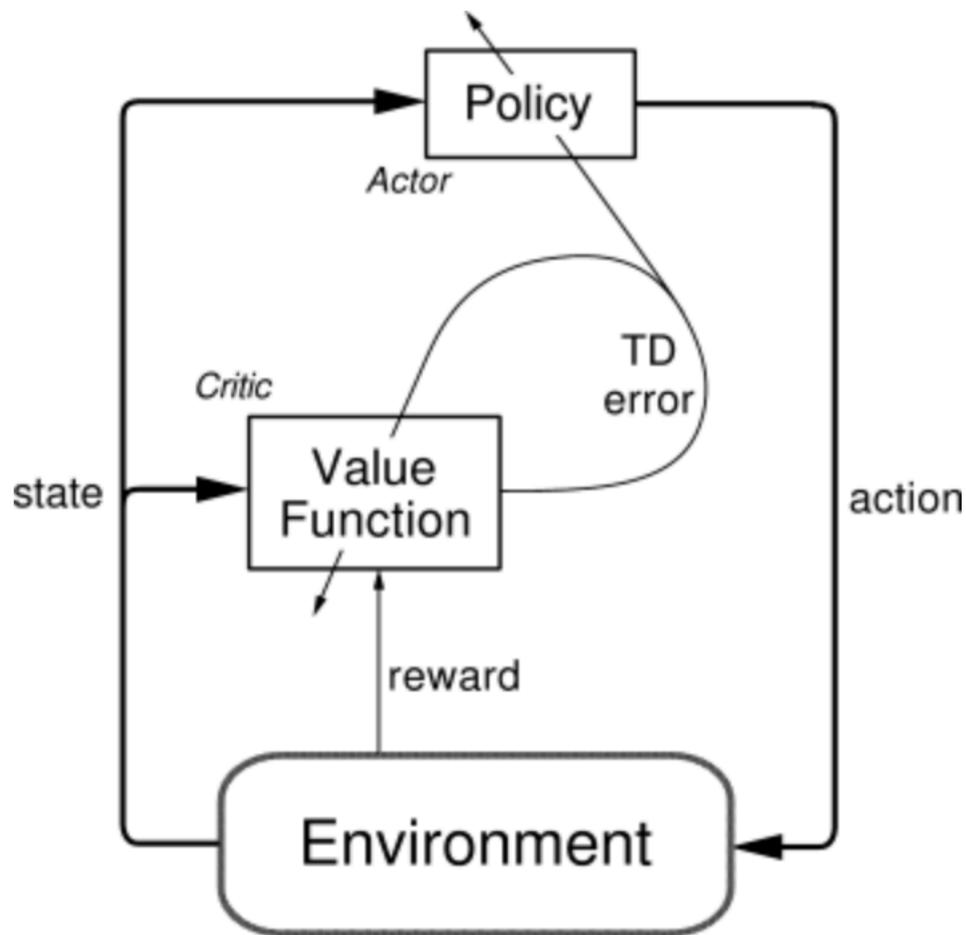
David Pfau, Oriol Vinyals. *Generative Adversarial Networks and Actor-Critic Methods.*

[arXiv:1610.01945](https://arxiv.org/abs/1610.01945)

Videos

Degrís – [Pendulum with pdf](#)

[Local Minimum Convergence](#)



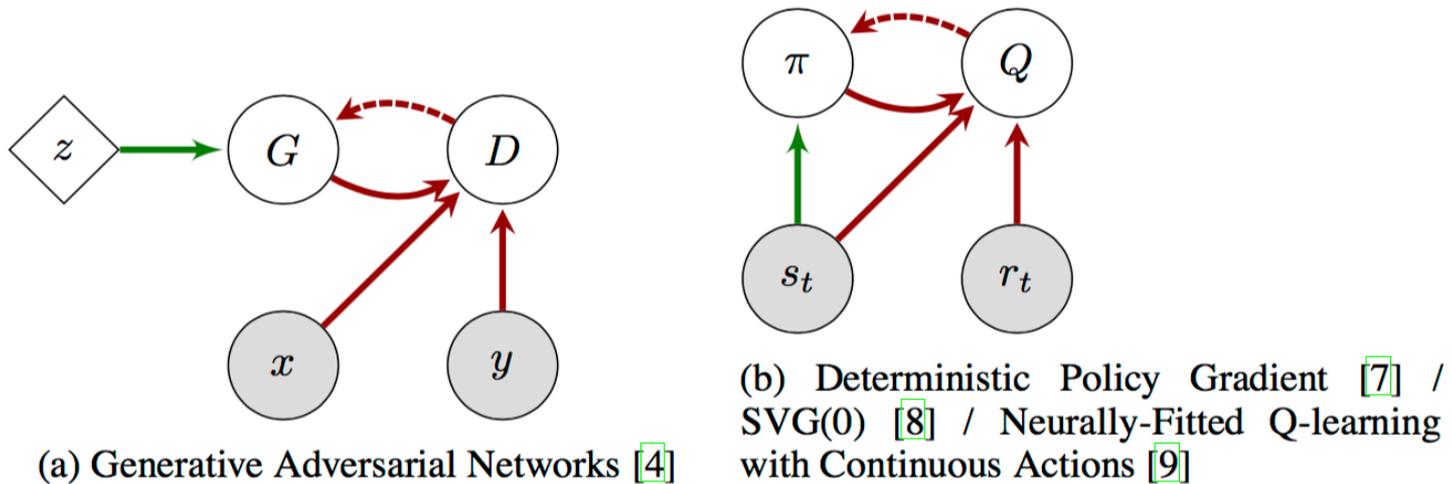


Figure 1: Information structure of GANs and AC methods. Empty circles represent models with a distinct loss function. Filled circles represent information from the environment. Diamonds represent fixed functions, both deterministic and stochastic. Solid lines represent the flow of information while dotted lines represent the flow of gradients used by another model. Paths which are analogous between the two models are highlighted in red. The dependence of Q on future states and the dependence of future states on π are omitted for clarity.

Both GANs and AC can be seen as bilevel or two-time-scale optimization problems, where one model is optimized with respect to the optimum of another model:

$$x^* = \arg \min_{x \in \mathcal{X}} F(x, y^*(x)) \quad (1)$$

$$y^*(x) = \arg \min_{y \in \mathcal{Y}} f(x, y) \quad (2)$$

$$Q^\pi(s, a) = \mathbb{E}_{s_{t+k} \sim \mathcal{P}, r_{t+k} \sim \mathcal{R}, a_{t+k} \sim \pi} \left[\sum_{k=1}^{\infty} \gamma^k r_{t+k} \mid s_t = s, a_t = a \right] \quad (6)$$

and learn a policy that is optimal for that value function:

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{s_0 \sim p_0, a_0 \sim \pi} [Q^\pi(s_0, a_0)] \quad (7)$$

We can express Q^π as the solution to a minimization problem:

$$Q^\pi = \arg \min_Q \mathbb{E}_{s_t, a_t \sim \pi} [\mathcal{D}(\mathbb{E}_{s_{t+1}, r_t, a_{t+1}} [r_t + \gamma Q(s_{t+1}, a_{t+1})] \parallel Q(s_t, a_t))] \quad (8)$$

Where $\mathcal{D}(\cdot \parallel \cdot)$ is any divergence that is positive except when the two are equal. Now the actor-critic problem can be expressed as a bilevel optimization problem as well:

$$F(Q, \pi) = \mathbb{E}_{s_t, a_t \sim \pi} [\mathcal{D}(\mathbb{E}_{s_{t+1}, r_t, a_{t+1}} [r_t + \gamma Q(s_{t+1}, a_{t+1})] \parallel Q(s_t, a_t))] \quad (9)$$

$$f(Q, \pi) = -\mathbb{E}_{s_0 \sim p_0, a_0 \sim \pi} [Q^\pi(s_0, a_0)] \quad (10)$$