

23/03/23

Lecture 03

Actor critic methods in Deep RL

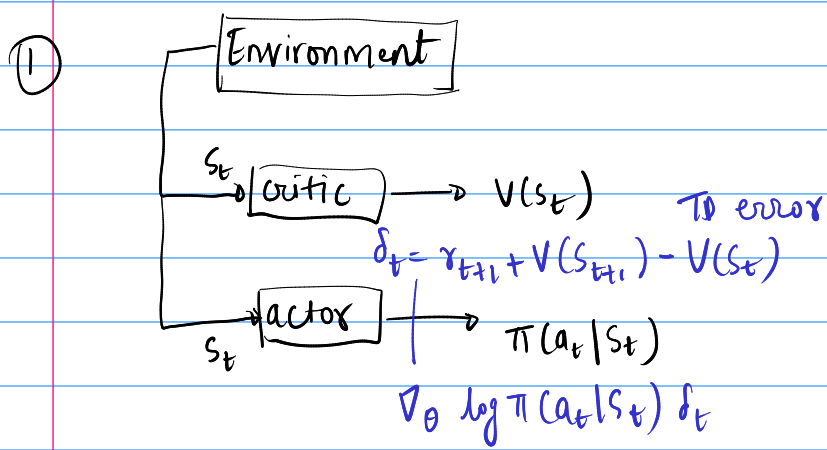
Basic policy gradient: $\nabla_{\theta} \log \pi(a_t | s_t; \theta) G_t$ ✓ called R_t in the A3C papers

Difficult to distinguish between G_t when results are similar, more stable gradients:

$$\nabla_{\theta} \log \pi(a_t | s_t; \theta) \underbrace{(G_t - b_{\theta}(s_t))}_{\text{Advantage}}$$

For us, $\theta \rightarrow$ parameters (weights, biases) of a neural network
 $\nabla_{\theta} \rightarrow$ gradient ascent

Two implementation options



* In both cases, one can use a common backbone for actor & critic

