

# Double-GAE

Jacob Hilton

October 15, 2020

We'd like to train autoregressive models using RL. But we have a choice: we can think of either completions or tokens as actions. This gives rise to two different policy gradient estimators. Here we'll analyze these estimators, and put those that use GAE [Schulman et al., 2015] into a unified perspective, using an advantage estimator we'll call Double-GAE.

## 1 Notation and terminology

Let's call completions *outer actions* and denote them using boldface letters, and let's call tokens *inner actions* and denote them using lightface letters with two indices, a completion index and a token index. Likewise for states, rewards and so on. We ignore the distinction between states and observations, and so states are prompts. We write states and actions next to one another to indicate concatenation. Thus

$$\mathbf{a}_t = a_{t0}a_{t1} \dots a_{t(n_t-1)},$$

where  $n_t$  is the length of the  $t$ th completion, and

$$s_{tu} = \mathbf{s}_t a_{t0} a_{t1} \dots a_{t(u-1)}$$

for  $0 \leq u \leq n_t - 1$ . We assume that the outer and inner rewards are related by

$$\mathbf{r}_t = r_{t0} + r_{t1} + \dots + r_{t(n_t-1)},$$

and we refer to  $r_{tu}$  for  $u < n_t - 1$  as *inner-only rewards*.

## 2 Double-GAE

Suppose our inner policy is given by an autoregressive model  $\pi_\theta(a_{tu} | s_{tu})$ . Then the outer policy is given by

$$\pi_\theta(\mathbf{a}_t | \mathbf{s}_t) = \pi_\theta(a_{t0} | s_{t0}) \pi_\theta(a_{t1} | s_{t1}) \dots \pi_\theta(a_{t(n_t-1)} | s_{t(n_t-1)}).$$

It follows that the outer policy gradient estimator

$$\hat{\mathbb{E}}_t \left[ \nabla_\theta \log \pi_\theta(\mathbf{a}_t | \mathbf{s}_t) \hat{\mathbf{A}}_t \right] = \hat{\mathbb{E}}_t \left[ \sum_{u=0}^{n_t-1} \nabla_\theta \log \pi_\theta(a_{tu} | s_{tu}) \hat{\mathbf{A}}_t \right] = \hat{\mathbb{E}}_t [r_t] \hat{\mathbb{E}}_{tu} \left[ \nabla_\theta \log \pi_\theta(a_{tu} | s_{tu}) \hat{\mathbf{A}}_t \right],$$

where  $\hat{\mathbf{A}}_t$  is the outer advantage estimator. Hence the outer policy gradient estimator can be viewed (up to a global constant) as a special case of the inner policy gradient estimator by using the inner advantage estimator  $\hat{A}_{tu} = \hat{\mathbf{A}}_t$  for all  $u$ . So there is no loss of generality by considering only the inner policy gradient estimator, as long as we allow non-standard inner advantage estimators.

Two natural choices for the inner advantage estimator are standard GAE, which we will call *inner GAE*, and the estimator obtained by using GAE with the outer policy gradient estimator, which we will call *outer GAE*. We would like to view both of these as special cases of a more general advantage estimator. To do this,

we introduce Double-GAE  $(\gamma_{\text{in}}, \gamma_{\text{out}}, \lambda_{\text{in}}, \lambda_{\text{out}})$ , which is parameterized by inner and outer discount rates  $\gamma_{\text{in}}$  and  $\gamma_{\text{out}}$  and bootstrapping parameters  $\lambda_{\text{in}}$  and  $\lambda_{\text{out}}$ . First we define the one-step backups

$$\delta_{tu}^V := \begin{cases} -V(s_{tu}) + r_{tu} + \gamma_{\text{in}} V(s_{t(u+1)}), & u < n_t - 1, \\ -V(s_{tu}) + r_{tu} + \gamma_{\text{out}} \gamma_{\text{in}}^{-(n_t-1)} V(s_{(t+1)0}), & u = n_t - 1. \end{cases}$$

Using these, we define the Double-GAE advantage estimator by

$$\hat{A}_{tu} := \sum_{k=0}^{\infty} \sum_{v=*}^{n_t+k-1} (\gamma_{\text{out}} \lambda_{\text{out}})^k (\gamma_{\text{in}} \lambda_{\text{in}})^{v-u} \delta_{(t+k)v}^V, \quad \text{where} \quad * := \begin{cases} u, & k = 0 \\ 0, & k > 0. \end{cases}$$

To recover inner GAE, we must assume that  $n_t = n$  is constant. If moreover  $\gamma_{\text{out}} = \gamma_{\text{in}}^n$  and  $\lambda_{\text{out}} = \lambda_{\text{in}}^n$ , then  $\delta_{tu}^V$  reduces to the standard inner one-step backup, and  $\hat{A}_{tu}$  reduces to the standard inner GAE estimator. Hence Double-GAE  $(\gamma_{\text{in}}, \gamma_{\text{in}}^n, \lambda_{\text{in}}, \lambda_{\text{in}}^n)$  is equivalent to inner GAE  $(\gamma_{\text{in}}, \lambda_{\text{in}})$ .

To recover outer GAE, we must assume both that there are no inner-only rewards, and that we only have an outer value function  $\mathbf{V}$ , from which we obtain our inner value function by taking  $V(s_{tu}) = \mathbf{V}(\mathbf{s}_t)$ . If moreover  $\gamma_{\text{in}} = \lambda_{\text{in}} = 1$ , then  $\delta_{tu}^V = 0$  for  $u < n_t - 1$  while  $\delta_{t(n_t-1)}^V$  reduces to the standard outer one-step backup, and  $\hat{A}_{tu}$  becomes independent of  $u$ , reducing to the standard outer GAE estimator for step  $t$ . It follows that inner Double-GAE  $(1, \gamma_{\text{out}}, 1, \lambda_{\text{out}})$  is equivalent (up to a global constant) to outer GAE  $(\gamma_{\text{out}}, \lambda_{\text{out}})$ .

We conclude from this that Double-GAE is general enough to describe most inner advantage estimators we are likely to care about, while keeping the number of hyperparameters minimal. Even though in practice  $n_t$  may not be constant, and we may consider it unhelpful to use a value function that can only differentiate between outer states, we have hopefully captured the most important features of any useful estimators.

### 3 Simplifying Double-GAE

We can split the double sum in the definition of the Double-GAE advantage estimator into a single sum and a double sum whose inner terms do not depend on  $u$ :

$$\begin{aligned} \hat{A}_{tu} &= \sum_{l=0}^{n_t-1-u} (\gamma_{\text{in}} \lambda_{\text{in}})^l \delta_{t(u+l)}^V + (\gamma_{\text{in}} \lambda_{\text{in}})^{-u} \sum_{k=1}^{\infty} (\gamma_{\text{out}} \lambda_{\text{out}})^k \sum_{v=0}^{n_t+k-1} (\gamma_{\text{in}} \lambda_{\text{in}})^v \delta_{(t+k)v}^V \\ &= \hat{A}_{tu}^{\text{comp}} + (\gamma_{\text{in}} \lambda_{\text{in}})^{-u} \sum_{k=1}^{\infty} (\gamma_{\text{out}} \lambda_{\text{out}})^k \hat{A}_{(t+k)0}^{\text{comp}}, \end{aligned}$$

where  $\hat{A}_{tu}^{\text{comp}}$  is defined by

$$\hat{A}_{tu}^{\text{comp}} := \sum_{l=0}^{n_t-1-u} (\gamma_{\text{in}} \lambda_{\text{in}})^l \delta_{t(u+l)}^V.$$

$\hat{A}_{tu}^{\text{comp}}$  is similar to the standard GAE estimator truncated after the end of the current completion, the only difference being the presence of  $\gamma_{\text{out}}$  in  $\delta_{t(n_t-1)}^V$ .

If  $\gamma_{\text{in}} = \lambda_{\text{in}} = 1$ , then  $\hat{A}_{tu}^{\text{comp}}$  simplifies to

$$\hat{A}_{tu}^{\text{comp}} = -V(s_{tu}) + r_{tu} + r_{t(u+1)} + \cdots + r_{t(n_t-1)} + \gamma_{\text{out}} V(s_{(t+1)0}).$$

In particular,  $\hat{A}_{t0}^{\text{comp}}$  is simply the standard outer one-step backup,

$$\hat{A}_{t0}^{\text{comp}} = \delta_t^V := -V(\mathbf{s}_t) + \mathbf{r}_t + \gamma_{\text{out}} V(\mathbf{s}_{t+1}).$$

Hence  $\hat{A}_{tu}$  simplifies to

$$\hat{A}_{tu} = -[-V(s_{t0}) + r_{t0} + r_{t1} + \cdots + r_{t(u-1)} + V(s_{tu})] + \sum_{k=0}^{\infty} (\gamma_{\text{out}} \lambda_{\text{out}})^k \delta_{t+k}^V.$$

This is almost identical to standard outer GAE, differing only by the  $u$ -step backup correction term in square brackets.

## 4 Choice of hyperparameters

It would be perverse for  $\gamma_{\text{in}}^{n_t}$  to be much smaller than  $\gamma_{\text{out}}$ , since that would involve putting more weight on rewards further into the future. Furthermore, if completions are not particularly long, then taking  $\gamma_{\text{in}}$  to be slightly less than 1 will not provide much variance reduction anyway. This suggests that taking  $\gamma_{\text{in}} = 1$  is a safe choice.

One could argue similarly for taking  $\lambda_{\text{in}} = 1$ , which would allow us to use the simplified calculation. However,  $\lambda$  is typically much lower than  $\gamma$ , so this may be a poor choice when completions are longer than a few tens of tokens, say.

It is hard to say much more than this without empirical data, but our overall recommendations are:

- For completions of just a few tokens, take  $\gamma_{\text{in}} = \lambda_{\text{in}} = 1$  (enabling the simplified calculation), and use reasonable defaults for  $\gamma_{\text{out}}$  and  $\lambda_{\text{out}}$  such as  $\gamma_{\text{out}} = 0.999$ ,  $\lambda_{\text{out}} = 0.95$  or a little lower.
- For completions much longer than a few tens of tokens, consider significantly lowering  $\lambda_{\text{out}}$  and taking  $\lambda_{\text{in}} = \lambda_{\text{out}}^{1/n}$ , where  $n$  is the number of tokens in a typical completion. It may also help to lower  $\gamma_{\text{out}}$ .
- For completions of an intermediate length, the first option might work fine, perhaps with lowering  $\lambda_{\text{out}}$ , but the second option could be worth trying.

## References

J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015.