

plain_text (0.95)

The generalized advantage estimator $\text{GAE}(\gamma, \lambda)$ is defined as the exponentially-weighted average of these k -step estimators:

$$\begin{aligned}\hat{A}_t^{\text{GAE}(\gamma, \lambda)} &:= (1 - \lambda) \left(\hat{A}_t^{(1)} + \lambda \hat{A}_t^{(2)} + \lambda^2 \hat{A}_t^{(3)} + \dots \right) \\ &= (1 - \lambda) \left(\delta_t^V + \lambda(\delta_t^V + \gamma \delta_{t+1}^V) + \lambda^2(\delta_t^V + \gamma \delta_{t+1}^V + \gamma^2 \delta_{t+2}^V) + \dots \right) \\ &= (1 - \lambda) \left(\delta_t^V (1 + \lambda + \lambda^2 + \dots) + \gamma \delta_{t+1}^V (\lambda + \lambda^2 + \lambda^3 + \dots) \right. \\ &\quad \left. + \gamma^2 \delta_{t+2}^V (\lambda^2 + \lambda^3 + \lambda^4 + \dots) + \dots \right) \\ &= (1 - \lambda) \left(\delta_t^V \left(\frac{1}{1 - \lambda} \right) + \gamma \delta_{t+1}^V \left(\frac{\lambda}{1 - \lambda} \right) + \gamma^2 \delta_{t+2}^V \left(\frac{\lambda^2}{1 - \lambda} \right) + \dots \right) \\ &= \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V\end{aligned}\tag{16}$$

plain_text (0.98)

From Equation (16), we see that the advantage estimator has a remarkably simple formula involving a discounted sum of Bellman residual terms. Section 4 discusses an interpretation of this formula as the returns in an MDP with a modified reward function. The construction we used above is closely analogous to the one used to define $\text{TD}(\lambda)$ (Sutton & Barto, 1998), however $\text{TD}(\lambda)$ is an estimator of the value function, whereas here we are estimating the advantage function.

plain_text (0.94)

There are two notable special cases of this formula, obtained by setting $\lambda = 0$ and $\lambda = 1$:

$$\text{GAE}(\gamma, 0) : \hat{A}_t := \delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)\tag{17}$$

$$\text{GAE}(\gamma, 1) : \hat{A}_t := \sum_{l=0}^{\infty} \gamma^l \delta_{t+l} = \sum_{l=0}^{\infty} \gamma^l r_{t+l} - V(s_t)\tag{18}$$

plain_text (0.98)

$\text{GAE}(\gamma, 1)$ is γ -just regardless of the accuracy of V , but it has high variance due to the sum of terms. $\text{GAE}(\gamma, 0)$ is γ -just for $V = V^{\pi, \gamma}$ and otherwise induces bias, but it typically has much lower variance. The generalized advantage estimator for $0 < \lambda < 1$ makes a compromise between bias and variance, controlled by parameter λ .

plain_text (0.98)

We’ve described an advantage estimator with two separate parameters γ and λ , both of which contribute to the bias-variance tradeoff when using an approximate value function. However, they serve different purposes and work best with different ranges of values. γ most importantly determines the scale of the value function $V^{\pi, \gamma}$, which does not depend on λ . Taking $\gamma < 1$ introduces bias into the policy gradient estimate, regardless of the value function’s accuracy. On the other hand, $\lambda < 1$ introduces bias only when the value function is inaccurate. Empirically, we find that the best value of λ is much lower than the best value of γ , likely because λ introduces far less bias than γ for a reasonably accurate value function.

plain_text (0.95)

Using the generalized advantage estimator, we can construct a biased estimator of g^γ , the discounted policy gradient from Equation (6):

$$g^\gamma \approx \mathbb{E} \left[\sum_{l=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \hat{A}_t^{\text{GAE}(\gamma, \lambda)} \right] = \mathbb{E} \left[\sum_{l=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}^V \right],\tag{19}$$

plain_text (0.93)

where equality holds when $\lambda = 1$.

title (0.93)

4 INTERPRETATION AS REWARD SHAPING

plain_text (0.98)

In this section, we discuss how one can interpret λ as an extra discount factor applied after forming a reward shaping transformation on the MDP. We also introduce the notion of a response function to help understand the bias introduced by γ and λ .

plain_text (0.97)

Reward shaping (Ng et al., 1999) refers to the following transformation of the reward function on an MDP: let $\Phi : \mathcal{S} \rightarrow \mathbb{R}$ be an arbitrary scalar-valued function on state space, and define the transformed reward function \tilde{r} by

$$\tilde{r}(s, a, s') = r(s, a, s') + \gamma \Phi(s') - \Phi(s)\tag{20}$$