Here, the subscript of $\mathbb{E}$ enumerates the variables being integrated over, where states and actions are sampled sequentially from the dynamics model $P(s_{t+1} \mid s_t, a_t)$ and policy $\pi(a_t \mid s_t)$, respectively. The colon notation $a : b$ refers to the inclusive range $(a, a + 1, \ldots, b)$. These formulas are well known and straightforward to obtain; they follow directly from Proposition 1, which will be stated shortly.

The choice $\Psi_t = A^\pi(s_t, a_t)$ yields almost the lowest possible variance, though in practice, the advantage function is not known and must be estimated. This statement can be intuitively justified by the following interpretation of the policy gradient: that a step in the policy gradient direction should increase the probability of better-than-average actions and decrease the probability of worse-than-average actions. The advantage function, by it's definition $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$, measures whether or not the action is better or worse than the policy's default behavior. Hence, we should choose $\Psi_t$ to be the advantage function $A^\pi(s_t, a_t)$, so that the gradient term $\Psi_t \nabla_\theta \log \pi_\theta(a_t \mid s_t)$ points in the direction of increased $\pi_\theta(a_t \mid s_t)$ if and only if $A^\pi(s_t, a_t) > 0$. See Greensmith et al. (2004) for a more rigorous analysis of the variance of policy gradient estimators and the effect of using a baseline.

We will introduce a parameter $\gamma$ that allows us to reduce variance by downweighting rewards corresponding to delayed effects, at the cost of introducing bias. This parameter corresponds to the discount factor used in discounted formulations of MDPs, but we treat it as a variance reduction parameter in an undiscounted problem; this technique was analyzed theoretically by Marbach & Tsitsiklis (2003); Kakade (2001b); Thomas (2014). The discounted value functions are given by:

$$V^{\pi,\gamma}(s_t) := \mathbb{E}_{\substack{s_{t+1:\infty}, \\ a_{t:\infty}}}\left[\sum_{l=0}^{\infty} \gamma^l r_{t+l}\right] \qquad Q^{\pi,\gamma}(s_t, a_t) := \mathbb{E}_{\substack{s_{t+1:\infty}, \\ a_{t+1:\infty}}}\left[\sum_{l=0}^{\infty} \gamma^l r_{t+l}\right] \qquad (4)$$

$$A^{\pi,\gamma}(s_t, a_t) := Q^{\pi,\gamma}(s_t, a_t) - V^{\pi,\gamma}(s_t) \qquad (5)$$

The discounted approximation to the policy gradient is defined as follows:

$$g^\gamma := \mathbb{E}_{\substack{s_{0:\infty} \\ a_{0:\infty}}}\left[\sum_{t=0}^{\infty} A^{\pi,\gamma}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t \mid s_t)\right]. \qquad (6)$$

The following section discusses how to obtain biased (but not too biased) estimators for $A^{\pi,\gamma}$, giving us noisy estimates of the discounted policy gradient in Equation (6).

Before proceeding, we will introduce the notion of a $\gamma$-just estimator of the advantage function, which is an estimator that does not introduce bias when we use it in place of $A^{\pi,\gamma}$ (which is not known and must be estimated) in Equation (6) to estimate $g^\gamma$.[1] Consider an advantage estimator $\hat{A}_t(s_{0:\infty}, a_{0:\infty})$, which may in general be a function of the entire trajectory.

**Definition 1.** *The estimator $\hat{A}_t$ is $\gamma$-just if*

$$\mathbb{E}_{\substack{s_{0:\infty} \\ a_{0:\infty}}}\left[\hat{A}_t(s_{0:\infty}, a_{0:\infty}) \nabla_\theta \log \pi_\theta(a_t \mid s_t)\right] = \mathbb{E}_{\substack{s_{0:\infty} \\ a_{0:\infty}}}\left[A^{\pi,\gamma}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t \mid s_t)\right]. \qquad (7)$$

It follows immediately that if $\hat{A}_t$ is $\gamma$-just for all $t$, then

$$\mathbb{E}_{\substack{s_{0:\infty} \\ a_{0:\infty}}}\left[\sum_{t=0}^{\infty} \hat{A}_t(s_{0:\infty}, a_{0:\infty}) \nabla_\theta \log \pi_\theta(a_t \mid s_t)\right] = g^\gamma \qquad (8)$$

One sufficient condition for $\hat{A}_t$ to be $\gamma$-just is that $\hat{A}_t$ decomposes as the difference between two functions $Q_t$ and $b_t$, where $Q_t$ can depend on any trajectory variables but gives an unbiased estimator of the $\gamma$-discounted $Q$-function, and $b_t$ is an arbitrary function of the states and actions sampled before $a_t$.

**Proposition 1.** *Suppose that $\hat{A}_t$ can be written in the form $\hat{A}_t(s_{0:\infty}, a_{0:\infty}) = Q_t(s_{t:\infty}, a_{t:\infty}) - b_t(s_{0:t}, a_{0:t-1})$ such that for all $(s_t, a_t)$, $\mathbb{E}_{s_{t+1:\infty}, a_{t+1:\infty} \mid s_t, a_t}[Q_t(s_{t:\infty}, a_{t:\infty})] = Q^{\pi,\gamma}(s_t, a_t)$. Then $\hat{A}$ is $\gamma$-just.*

---

[1]Note, that we have already introduced bias by using $A^{\pi,\gamma}$ in place of $A^\pi$; here we are concerned with obtaining an unbiased estimate of $g^\gamma$, which is a biased estimate of the policy gradient of the undiscounted MDP.