**Algorithm 1** Simplified Vanilla Policy Gradient Algorithm

---

1: Initialize the policy $\pi$ and the value function $V$.
2: **for** $k = 0, 1, 2, ...$ **do**
3:     Collect set of trajectories $D_k = \{\tau_i\}$ by running
                       policy $\pi_k$ in the environment.
4:     Compute future returns $\hat{R}_t$ for each time step $t$ in each trajectory $\tau_i$.
5:     Compute the policy gradient $\hat{g}_k$ using $D_k$ and $V$.    ▷ More on $\hat{g}_k$ soon!
6:     Update the policy $\pi$ with $\hat{g}_k$ using SGD, Adam, etc.
7:     Fit value function $V$ by regression on MSE between
                       $V$ and $\hat{R}_t$ using SGD, Adam, etc.
8: **end for**

---