

title (0.90)

## 6.1 POLICY OPTIMIZATION ALGORITHM 0

plain\_text (0.97)

While generalized advantage estimation can be used along with a variety of different policy gradient methods, for these experiments, we performed the policy updates using trust region policy optimization (TRPO) (Schulman et al., 2015). TRPO updates the policy by approximately solving the following constrained optimization problem each iteration:

$$\begin{aligned}
 & \underset{\theta}{\text{minimize}} \quad L_{\theta_{old}}(\theta) \\
 & \text{subject to} \quad \overline{D}_{\text{KL}}^{\theta_{old}}(\pi_{\theta_{old}}, \pi_{\theta}) \leq \epsilon \\
 & \text{where } L_{\theta_{old}}(\theta) = \frac{1}{N} \sum_{n=1}^N \frac{\pi_{\theta}(a_n | s_n)}{\pi_{\theta_{old}}(a_n | s_n)} \hat{A}_n \\
 & \quad \overline{D}_{\text{KL}}^{\theta_{old}}(\pi_{\theta_{old}}, \pi_{\theta}) = \frac{1}{N} \sum_{n=1}^N D_{\text{KL}}(\pi_{\theta_{old}}(\cdot | s_n) \| \pi_{\theta}(\cdot | s_n))
 \end{aligned} \tag{31}$$

plain\_text (0.98)

As described in (Schulman et al., 2015), we approximately solve this problem by linearizing the objective and quadraticizing the constraint, which yields a step in the direction  $\theta - \theta_{old} \propto -F^{-1}g$ , where  $F$  is the average Fisher information matrix, and  $g$  is a policy gradient estimate. This policy update yields the same step direction as the natural policy gradient (Kakade, 2001a) and natural actor-critic (Peters & Schaal, 2008), however it uses a different stepsize determination scheme and numerical procedure for computing the step.

plain\_text (0.97)

Since prior work (Schulman et al., 2015) compared TRPO to a variety of different policy optimization algorithms, we will not repeat these comparisons; rather, we will focus on varying the  $\gamma, \lambda$  parameters of policy gradient estimator while keeping the underlying algorithm fixed.

plain\_text (0.73)

For completeness, the whole algorithm for iteratively updating policy and value function is given below:

table (0.43)

```

Initialize policy parameter  $\theta_0$  and value function parameter  $\phi_0$ .
for  $i = 0, 1, 2, \dots$  do
    Simulate current policy  $\pi_{\theta_i}$  until  $N$  timesteps are obtained.
    Compute  $\delta_t^V$  at all timesteps  $t \in \{1, 2, \dots, N\}$ , using  $V = V_{\phi_i}$ .
    Compute  $\hat{A}_t = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l}^V$  at all timesteps.
    Compute  $\theta_{i+1}$  with TRPO update, Equation (31).
    Compute  $\phi_{i+1}$  with Equation (30).
end for

```

plain\_text (0.98)

Note that the policy update  $\theta_i \rightarrow \theta_{i+1}$  is performed using the value function  $V_{\phi_i}$  for advantage estimation, not  $V_{\phi_{i+1}}$ . Additional bias would have been introduced if we updated the value function first. To see this, consider the extreme case where we overfit the value function, and the Bellman residual  $r_t + \gamma V(s_{t+1}) - V(s_t)$  becomes zero at all timesteps—the policy gradient estimate would be zero.

title (0.91)

## 6.2 EXPERIMENTAL SET 8

plain\_text (0.98)

We evaluated our approach on the classic cart-pole balancing problem, as well as several challenging 3D locomotion tasks: (1) bipedal locomotion; (2) quadrupedal locomotion; (3) dynamically standing up, for the biped, which starts off laying on its back. The models are shown in Figure 1.

title (0.91)

## 6.2.1 ARCHITECTURE 10

plain\_text (0.98)

We used the same neural network architecture for all of the 3D robot tasks, which was a feedforward network with three hidden layers, with 100, 50 and 25 tanh units respectively. The same architecture was used for the policy and value function. The final output layer had linear activation. The value function estimator used the same architecture, but with only one scalar output. For the simpler cart-pole task, we used a linear policy, and a neural network with one 20-unit hidden layer as the value function.