

算法概览图符号表

符号	英文含义	中文含义
t	timestep in real environment	在真实环境中的时间步
$\gamma \in [0, 1]$	discount factor	折扣因子
k	hypothetical step	“假设步” 结果在模型中
o_t	observation at timestep t	在时间步 t 的观察
u_t	reward at timestep t	在时间步 t 的奖励
a_t	action at timestep t	在时间步 t 的动作
s_t	hidden state at timestep t	在时间步 t 的隐藏状态
r_t	predicted reward at timestep t	在时间步 t 的预测奖励
p	predicted policy logits	预测的策略对数几率
v	predicted value	预测的值
h	Representation Network	表征网络, hidden state
g	Dynamics Network	动力学网络
f	Prediction Network	预测网络
vp	predicted value prefix	预测的值前缀
r_h^k	Reward hidden state at hypothetical step k	在假设步 k 的奖励隐藏状态 Efficient2
g^{rnn}	(recurrent) Dynamics Network	(循环) 动力学网络 vp^k, s^k, a^k Efficient2
β	Proposal policy distribution	提议策略分布。
$\hat{\beta}$	empirical policy distribution	经验策略分布

$\tau = \{s_0, a_0, s_1, a_1, \dots\}$	trajectory	
$R(\tau) = \sum_{t=1}^T r_t$	return	
$G_t = \sum_{k=0}^{T-t-1} \gamma^k r_{t+k}$	Discounted return	
$\pi\theta(a_t s_t)$	stochastic policy	随机策略
$V_{\pi}(s) = \mathbb{E}_{\pi}[G_t s_t = s]$	state value fuction	
$Q_{\pi}(s, a) = \mathbb{E}_{\pi}[G_t s_t = s, a_t = a]$	state-action value fuction	

(表1：算法概览图的符号说明。)