

# REINFORCE & Baseline - REINFORCE

Policy fn approximator: Neural Net with 1 hidden layer

$$W_{in}: n\_hidden \times n\_state\_dims$$

$$b_{in}: n\_hidden$$

$$W_{out}: n\_actions \times n\_hidden$$

$$b_{out}: n\_actions$$

} parameters

let  $s_t: n\_state\_dims \in S$

$$\pi(s_t) := \text{SOFTMAX}(W_{out} \cdot \text{RELU}(W_{in} \cdot s_t + b_{in}) + b_{out})$$

Prob. dist over  $\{1 \dots n\_actions\}$

\* there's also dropout of  $p = 0.6$  b/w the layers

$$n\_hidden = 128$$

$$\lambda: 0.01$$

$$\gamma \in [1, 0.99, 0.95]$$

Value fn approximator: NN with 1 hidden layer

$W_{in} : n\_hidden \times n\_state\_dims$

$b_{in} : n\_hidden$

$W_{out} : 1 \times n\_hidden$

$b_{out} : 1$

let  $s_t : n\_state\_dims \in S$

} parameters

$$V(s_t) := W_{out} \cdot \text{RELU}(W_{in} \cdot s_t + b_{in}) + b_{out}$$

\* there's also dropout of  $p = 0.6$  b/w the layers

\*  $n\_hidden = 128$

$l_n : 0.01$

$\gamma \in [1, 0.99, 0.95]$