

The objective function:

$$\min_{\theta} \sum_{t=1}^N R_t \left(- \sum_{i=1}^n \mathbf{I}(a_i, a_t) \log(\pi(a_i | s_t; \theta)) \right)$$

where N is the number of time step in one episode, n is the number of actions in action space, $\mathbf{I}(x, y)$ is an indicator function defined as $\mathbf{I}(x, y) = 1$ if $x == y$ else 0.