Learner

Objective

$$\pi(s) = \mathtt{softmax}(g_{\theta}(s))$$

$$\mathcal{L}(a, \pi(s)) = H(a, \pi(s))$$

Hypothesis space

Convolutional neural net

Optimizer

Stochastic gradient descent

$$\rightarrow \pi: s \rightarrow a$$

$$\pi^* = \underset{\pi \in \Pi}{\operatorname{arg\,min}} \sum_{i=1}^{N} \mathcal{L}(\pi(s_i), a_i)$$