

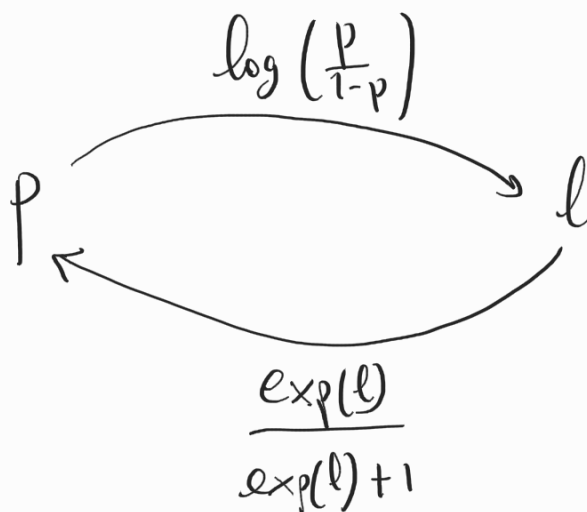
# ① Function approximation

probability  
↓  
logit

$p$   
↓  
 $\log\left(\frac{p}{1-p}\right)$

$$(0, 1)$$

$$(-\infty, +\infty)$$

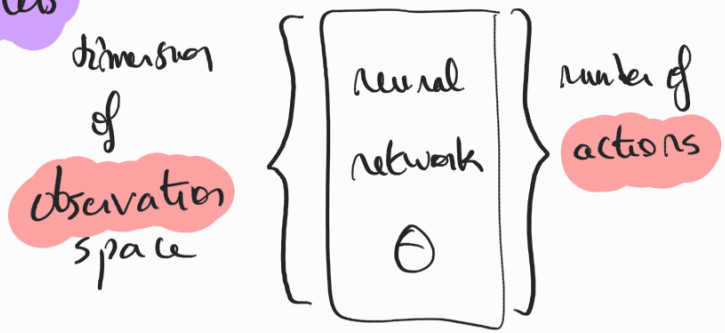
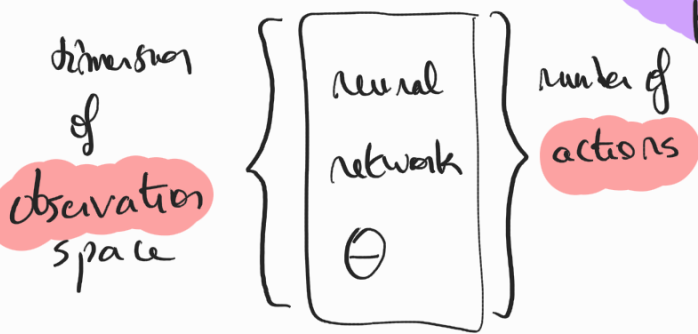


Question: What if  $S$  is infinite?

↳ change of representation

Deterministic policy

Stochastic policy



$$q_{\theta}: S \times A \rightarrow \mathbb{R} \quad \leftarrow \text{value}$$

$$\sigma_{\theta}(s) = \arg \max_{a \in A} q(s, a)$$

$$p_{\theta}: S \times A \rightarrow \mathbb{R} \quad \leftarrow \text{logit}$$

$$\sigma_{\theta}(s) = \text{distribution}(a \mapsto p(s, a))$$

apply softmax here to get probabilities

Update: parameters

$$[ \text{NEW} = \text{OLD} + \alpha (\text{CURRENT} - \text{OLD}) ]$$

↳ includes "gradient ascent algorithms"

Formulation:

$$\min_{\theta} \mathcal{L}(\theta)$$

loss

Typical approach: stochastic gradient descent

Sketch:

(1) Batch sampling:

using the current policy, find SETS of trajectories

(2) Batch update:

(4) Batch update.

using the batch, update the parameters