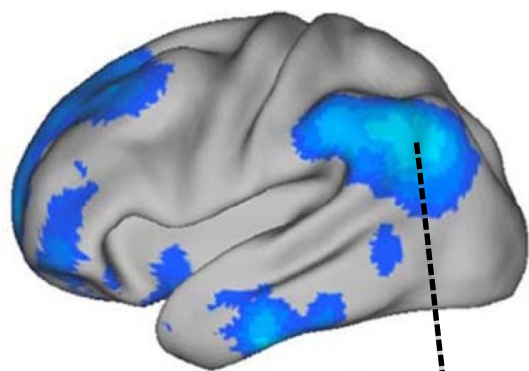


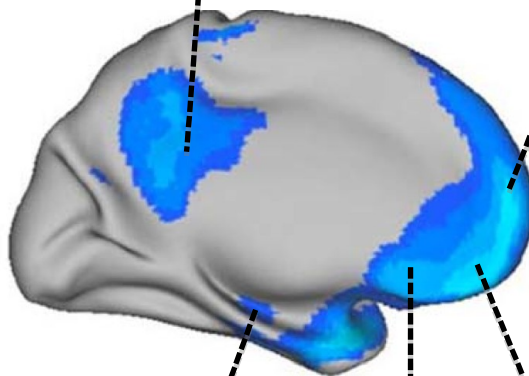
$$\nabla_{\theta_{k+1}^Q} \mathcal{L}(\theta_{k+1}^Q) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{U}(\mathcal{D})} \left[\underbrace{(Q(s,a|\theta_{k+1}^Q) - y)}_{\text{regret}} \underbrace{\nabla_{\theta_{k+1}^Q} Q(s,a|\theta_{k+1}^Q)}_{???} \right]$$

parameter update (backprop!)



s, a (really ?)

world knowledge



r_t

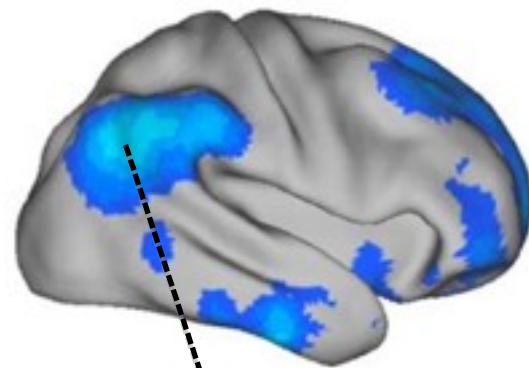
Instantaneous reward

$(s, a, r, s') \sim \mathcal{U}(\mathcal{D})$

experience sampling

$Q(s, a|\theta_k^Q)$

value function



$a = \operatorname{argmax}_{\tilde{a} \in \mathcal{A}} Q(s, \tilde{a}|\theta_k^Q) + u(s, \tilde{a})$

policy function (greedy!)