reinforcement learning (learning what to do)

"Life is like playing a violin solo in public while learning the instrument as one goes on"



① understand challenges of decision-making

2 learn basics of RL framework

③ see applications to moral psychology

in the beginning...



N





the moral universe (1/2)



Altruistic punishment in humans

Ernst Fehr* & Simon Gächter†

* University of Zürich, Institute for Empirical Research in Economics, Blümlisalpstrasse 10, CH-8006 Zürich, Switzerland † University of St Gallen, FEW-HSG, Varnbüelstrasse 14, CH-9000 St Gallen, Switzerland

REVIEW ARTICLE

articles

Punishment in animal societies

T. H. Clutton-Brock & G. A. Parker



learning + deciding





punisher



Expected Value = Reward × Probability $\Theta = p(P|T)$ $Q_T = B + C \times p(P|T)$ $Q_A = 0$ $p(T) = \frac{e^{Q_T}}{\Sigma e^Q}$ p(T) $Q_T - Q_A$







PROBLEMS

① an awful lot of mental accounting to do





2 a fear of god or a taste for justice?





the reinforcement learning solution

prediction error learning
(c.f. rescorla-wagner)

 $PE = R - Q_t \qquad R = B - C$

 $Q_t \leftarrow Q_t + \alpha(PE)$



history →

PROBLEMS

① an awful lot of mental accounting to do





② a fear of god or a taste for justice?



Reinforcement Learning



the reinforcement learning solution

prediction error learning
(c.f. rescorla-wagner)

$$PE = R - Q_{t-1} \qquad R = B - C$$

 $Q_t \leftarrow Q_t + \alpha(PE)$





allow the reward function to evolve

theft/punishment

history \rightarrow

so what happens?

① Resolute punishment (innately rewarding)

② Flexible theft ("objective" rewards only)

Temporal Difference Learning



Model Based Reinforcement Learning

L = + R = 0 L = + R = -

R = **0**

Model-Free

Temporal Credit Assignment Problem Solution: Treat good options like rewards $Q_{(s,a)} \leftarrow Q_{(s,a)} + \alpha PE$ $Q_{(s,a)} \leftarrow Q_{(s,a)} + \alpha (R_{s'} + Q_{(s',a')} - Q_{(s,a)})$ $0 = Q_{(\bullet,L)} \leftarrow Q_{(\bullet,L)} + \frac{1}{2}(R_{\bullet} + Q_{(\bullet,L)} - Q_{(\bullet,L)})$ $2 = Q_{(\bullet,L)} \leftarrow Q_{(\bullet,L)} + \frac{1}{2}(Q_{\bullet} + Q_{(\bullet,L)} - Q_{(\bullet,L)})$ $0 = Q_{(\bullet,L)} \leftarrow Q_{(\bullet,L)} + \frac{1}{2}(Q_{\bullet} + Q_{(\bullet,L)} - Q_{(\bullet,L)})$ $0 = Q_{(\bullet,L)} \leftarrow Q_{(\bullet,L)} + \frac{1}{2}(Q_{\bullet,L} + Q_{(\bullet,L)} - Q_{(\bullet,L)})$ $0 = Q_{(\bullet,L)} \leftarrow Q_{(\bullet,L)} + \frac{1}{2}(Q_{\bullet,L} + Q_{(\bullet,L)} - Q_{(\bullet,L)})$

Reinforcement Learning

From Association to Action

e.g. Dickinson et al 1995



Model-Free Reinforcement Learning



Model Based Reinforcement Learning



From Association to Action



Model-Free Reinforcement Learning



Model Based Reinforcement Learning



Prediction of value, but no model of outcome

Memory: Patient HM





Neural Mechanisms of Model-Free Learning

Graphic borrowed from Daw & Shohamy 2008 For experiments, see e.g. Shultz et al 1996



Dopamine = prediction error signal





Model-Free Learning





Nucleus Accumbens of the Striatum

Model-Free Learning



Nucleus Accumbens of the Striatum

Model-Free Learning





Nucleus Accumbens of the Striatum

Application: Why is cocaine addictiQe?

Redish, 2004





Olds & Milner 1964

Application: Why is cocaine addictiQe?

Redish, 2004









The AQersion to Harm



Wendy Mendes







Cushman, Gray, Gaffey & Mendes, 2012



Dual Process Morality

Josh Greene





