Action learning Instrumental conditioning

David G. Nagy - neural modelling - 2023













classical conditioning







- immediate choice
- evolutionary programming
- sequential decisions
- vigour



• immediate choice

- evolutionary programming
- sequential decisions
- vigour





0.5

P(L)





P(R)



P(L)







0.8

P(L)



P(R)

P(L) 0.8

P(L)



P(R) = 1 - P(L)



m

P(R) = 1 - P(L)

 $P(L \mid m)$



 $P(L \mid m)$



 $P(L \mid m)$











exploration

$$P(L) = \frac{e^{\beta m_L}}{e^{\beta m_R} + e^{\beta m_L}}$$





how should we set m?





 $m_L = \mathbb{E}[r \mid a = L] = \langle r_L \rangle$



update rule:

$$m_L \rightarrow m_L + \epsilon (r_L - m_L)$$

$$m_R \to m_R + \epsilon (r_R - m_R)$$









$\mathbb{E}[r | \mathbf{m}] = \langle r \rangle$ $= \mathbb{E}_{a \sim P(a | \mathbf{m})}[r_a]$ $= P(L) \langle r_L \rangle + P(R) \langle r_R \rangle$

update rule:

$$m_L \to m_L + \epsilon \Delta m_L$$
$$\mathbb{E}[\Delta m_L] \propto \frac{\partial \langle r \rangle}{\partial m_L}$$



$$\frac{\partial P(L)}{\partial m_L} = \beta P(L)(1 - P(L))$$

$$\frac{\partial P(R)}{\partial m_L} = -\beta P(R)P(L)$$

 $\frac{\partial \langle r \rangle}{dr} = \beta P(L)(\langle r_L \rangle - \langle r \rangle)$ ∂m_L





- immediate choice
- evolutionary programming
- sequential decisions
- vigour

The Skinnerian Pigeon







The Hapless Pigeon



Omission



active coping strategies evoked from the IPAG and the dIPAG













inhibition



invigoration





inhibition





invigoration

inhibition





- data
- inst = RW + noise
- inst + bias
- inst + bias + Pavlovian
 - go+ $\omega V(s_t)$

(Guitart-Masip et al, 2011; 2012)



instrumental control



Pavlovian control

(Dayan et al., 2006)



- immediate choice
- evolutionary programming
- sequential decisions
- vigour

The Problem of Time



Humboldt, Saskatchewan



'ethologically': maximize long term reward

A Simpler Example





requires a signal reporting the reward: dopamine

A Simpler Example





action reward









Two Ideas



• idea a: learn values V for states

r(<u>____</u>) = 4 is *immediately* valuable V() = 4 is *secondarily* valuable

Idea a: learn values V for states

- prediction: $V(\bullet)$









learning rate



Two Ideas



• idea a: learn values V for states

r(_____) = 4 is *immediately* valuable

V() = 4 is secondarily valuable

idea b: use values as surrogate feedback

- = 4 so do:
- →: V() = 3



Idea b: Use values as surrogate feedback



$$V(\bullet) = 4$$

$$V(t) = r(t)$$

$$\delta(t) = r(t) - V(t)$$
prediction error
$$V(\bullet) = 0 + V(\bullet)$$

$$V(t) = r(t) + V(t+1)$$

 $\delta(t) = r(t) + V(t+1) - V(t)$

Sutton: temporal difference prediction error Bellman: asynchronous dynamic programming



Prediction, Prediction Errors and Dopamine $\delta(t) = r(t) + V(t + 1) - V(t)$











Idea 3: Prediction for Control



- If: V(-) = 4 V(-) = 3
- V(**•**) = 3.5
 - $\delta(t) = r(t) + V(t+1) V(t)$
- $\delta(t) = 0 + 4 3.5 = +0.5$
- $\delta(t) = 0 + 3 3.5 = -0.5$

so increase the probability of **—**





actor-critic

critic

$m_b(u) \rightarrow (1 - \epsilon)m_b(u) + \epsilon_A \delta_{ab}\delta$ actor

 $m_b \rightarrow (1 - \epsilon)m_L + \epsilon \delta_{ab}(r_a - v)$



actor-critic





dopamine signals to both motivational & motor striatum appear, surprisingly the same

suggestion: training both values & policies

actor-critic

- direct actor
 - immediate choice version of policy gradient
- policy gradient for sequential decisions
 - REINFORCE algorithm, likelihood ratio PG, MC estimate
 - can be used without baseline, worse convergence
 - with baseline = V(s), we also have to learn value
 - learn V (via e.g. TD rule): critic
 - learn policy: actor

model based RL



model based RL



model based RL



- immediate choice
- evolutionary programming
- sequential decisions

• vigour

- Two components to choice: – what:
 - lever pressing
 - direction to run
 - meal to choose
 - when/how fast/how vigorous
 - free operant tasks
- real-valued dynamic programming

Vigour

The model

<u>Goal</u>: Choose actions and latencies to maximize the *average rate of return* (rewards minus costs per time)

Average Reward Cost/benefit Tradeoffs

1. Which action to take?

 \Rightarrow Choose action with largest expected reward minus cost

2. How fast to perform it?

slow → less costly (vigour cost)
 slow → delays (all) rewards

 \Rightarrow Choose rate that balances vigour and opportunity costs

explains faster (irrelevant) actions under hunger, etc

 net rate of rewards = cost of delay (opportunity cost of time)

Relation to Dopamine

Phasic dopamine firing = reward prediction error

What about tonic dopamine?

more

$$\sum \delta_t = \sum [r_{t+1} + V(s_{t+1}) - V(s_t)]$$

so the average is the same as the average total reward, the timing is just moved around

60

= $\sum [r_{t+1}] + V(s_T) - V(s_0)$...also explains effects of phasic dopamine on response times

+ direct evidence from Berke: pro and contra

Classical/Instrumental Conditioning

prediction: of important events control:

- Ethology
 - optimality
 - appropriateness
- Psychology
 - classical/operant conditioning
 - Neurobiology

neuromodulators; amygdala; OFC

- in the light of those predictions
 - Computation
 - dynamic progr.
 - Kalman filtering
 - Algorithm
 - TD/delta rules
 - simple weights
- nucleus accumbens; dorsal striatum

homework

- reproduce figures 9.8 (policy evaluation) and 9.9 (actor-critic) from the chapter for the environment of figure 9.7
- try actor critic for a maze task

homework

Tutorial by Noémi Éltető tomorrow at DZNE Lecture Room

