## Reinforcement Learning II: Biology

Quentin Huys

Wellcome Trust Centre for Neuroimaging, UCL Gatsby Computational Neuroscience Unit, UCL Guy's and St. Thomas' Hospital NHS Trust

Psychiatrische Universitatsklinik, University of Zurich Theoretical Neuromodeling Unit, ETH

Advanced Course in Computational Neuroscience, Bedlewo, Poland, August 2012

#### Recapitulation

- MDP: {s,a,T,R,pi}
- Multiple solution approaches
- Global": policy evaluation & iteration
- "Local": sampling:
  - Monte Carlo
  - TD
- How do humans or animals solve RL problems?

## Policy

- Unconditioned responses
- Pavlovian conditioning
- Habits
- Goal-directed behaviour
- Doing tree search



#### Unconditioned responses



#### are innate evolutionary strategies



#### Innate evolutionary strategies





Hirsch and Bolles 1980

#### Innate evolutionary strategies

## are quite sophisticated...



Hirsch and Bolles 1980

## Specific biological substrates





prior to training

CS: bell -> no response

US: food -> salivation & consumption



prior to training CS: bell -> no response US: food -> salivation & consumption training

CS: bell  $\longrightarrow$  US: food



prior to training CS: bell -> no response US: food -> salivation & consumption training

CS: bell  $\longrightarrow$  US: food





CS: bell -> salivation US: food -> salivation & consumption





prior to training CS: bell -> no response US: food -> salivation & consumption training CS: bell  $\longrightarrow$  US: food after training: CS: bell -> salivation US: food -> salivation & consumption omission training: salivation - US: omitted

no salivation



CS: bell

US: food

#### Aversive Pavlovian effects: freezing



#### Aversive Pavlovian effects: freezing



#### A UK checkout



#### A UK checkout





O Ministério da Saúde adverle: O uso deste produto obstrui as artérias e dificulta a circulação do sangue.



Reinforcement learning

#### A UK checkout

WARNING: Tobacco smoke can harm your children.





#### Tobacco firms accused of funding campaign to keep cigarettes on display

Jamie Doward and Alex Ascherson guardian.co.uk, Saturday 26 February 2011 20.54 GMT Article history



Cigarettes displayed for sale at a store in central London. Photograph: Andy Rain/EPA

The Guardian, Sat Feb 26th 2011

#### Pavlovian influences on instrumental learning



Guitart-Masip, Huys et al. 2011, 2012

Reinforcement learning

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

Quentin Huys

#### Pavlovian influences on instrumental learning



Guitart-Masip, Huys et al. 2011, 2012

#### Pavlovian influences on instrumental learning



#### Guitart-Masip, Huys et al. 2011, 2012

By being predictive of an affective outcome, a neutral stimulus can come to elicit the innate preparatory response usually evoked by the affective outcome.

p(active response) ~ value(stimulus)?

#### Model





#### Model



#### $p(\mathbf{go}|s_t) \propto \mathcal{Q}_t(\mathbf{go}|s_t) + \mathbf{bias}(\mathbf{go})$



#### Model



# $p(\mathbf{go}|s_t) \propto \mathcal{Q}_t(\mathbf{go}|s_t) + \mathbf{bias}(\mathbf{go}) + \mathcal{V}_t(s_t)$ $\mathcal{V}_t(s_t) = \mathcal{V}_{t-1}(s_t) + \epsilon(r_t - \mathcal{V}_{t-1}(s_t))$







#### Pavlovian values in the brain?

- Where the values V(s) are is not so clear
- Midbrain dopamine neurons seem to report a TD error



#### Pavlovian values in the brain?

- Where the values V(s) are is not so clear
- Midbrain dopamine neurons seem to report a TD error





#### Pavlovian values in the brain?

- Where the values V(s) are is not so clear
- Midbrain dopamine neurons seem to report a TD error















#### 'Sign-tracking'

#### 'Sign-tracking'

- Wiers et al. 2010: Attentional training reduces relapse
- Heinze et al. 2009: NAcc DBS reduces relapse & craving

#### Depression: aversive stimulus values

- dot probe task
- emotional stroop
- recollection bias
# Devaluation



# Goal-directed vs. habitual behaviour mix and match

# Habits

- Instrumental stimulus-response tendencies
- Lever pressing: not innate
- SARSA:

 $\mathcal{Q}(s_t, a_t) \leftarrow \mathcal{Q}(s_t, a_t) + \alpha[r_t + \gamma \mathcal{Q}(s_{t+1}, a_{t+1}) - \mathcal{Q}(s_t, a_t)]$ 

- any response
- slow acquisition
- slow 'extinction' if reward changes

#### Devaluation

- Train briefly / overtrain
- Devalue outcome
- Immediate behaviour?

early: g-d late: habits

#### Devaluation

- Train briefly / overtrain
- Devalue outcome
- Immediate behaviour?



#### Devaluation

- Train briefly / overtrain
- Devalue outcome
- Immediate behaviour?



#### Devaluation

- Train briefly / overtrain
- Devalue outcome
- Immediate behaviour?



## Simple is better at times: doctors



#### 20 cases for which truth known

Cardiologists General physicians A&E physicians

Melly et al. 2002

# Simple is better at times: doctors



#### 20 cases for which truth known

Cardiologists General physicians A&E physicians

Physicians overly cautious, but still miss many -> complications

Melly et al. 2002

- Also see phasic DA responses in instrumental learning
- DA seems to be involved in learning with prediction errors
  - Sum of prediction errors = "cached" value (Daw et al. 2005)

# Many decision systems in parallel



- behaviourally and neurobiologically distinct and identifiable
  - How to arbitrate? according to advantages and disadvantages?
- interactions
  - Pavlovian pruning of decision trees
  - model-based teaching of habits

# Goal-directed decision making



- 30<sup>40</sup>?
- Legal boards ~10<sup>123</sup>
- Can't just do full tree search.

01

a2

01 02 01 02



# Pruning a decision tree



## Pruning a decision tree



Don't go there... Don't think it either?

# Pavlovian pruning

- Optimality
  - conserve guarantees
  - difficult & computationally expensive

#### Approximate

 trade optimality for speed



# Pavlovian pruning

- Optimality
  - conserve guarantees
  - difficult & computationally expensive

#### Approximate

 trade optimality for speed



# Pavlovian

• reflexively prune on encountering a punishment

# **Psychochess**



# A poor experimental psychologist's version of chess



-X = -140

A tree search task

# A poor experimental psychologist's version of chess



-X = -140

A tree search task

# Model I: full lookahead















# People don't look to the end...



Huys et al. (2011) Submitted

**Reinforcement learning** 

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

Quentin Huys

## Chess



Selective pruning! Pavlovian pruning???





# Pavlovian?



# Pavlovian?







Reinforcement learning

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

Quentin Huys

# Pavlovian?







Reinforcement learning

Quentin Huys

# Stop at losses?



# Stop at losses?



# Stop at losses?



# **Avoiding losses**





Start state 3, 3 moves



# **Avoiding losses**





Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

# Adaptive pruning model


# Adaptive pruning model



#### Adaptive pruning wins



Huys et al. (2011) Submitted

**Reinforcement learning** 

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

#### Generate some data









 $\mathcal{V}(s) \leftarrow \mathcal{V}(s) + \epsilon(\mathcal{V}(s') + r_t - \mathcal{V}(s))$ 



 $\mathcal{V}(s) \leftarrow \mathcal{V}(s) + \epsilon(\mathcal{V}(s') + r_t - \mathcal{V}(s))$ 

#### Pavlovian influences inside a decision tree



# Adaptive pruning wins



Huys et al. (2011) Submitted

**Reinforcement learning** 

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

#### **Further alternatives**

#### Loss aversion

- maybe subjects dislike large loss disproportionately?
- Cached terminal values
- Cached learning only
- Choices of entire sequences
- Compare models' abilities to explain the ENTIRE set of data



Huys et al. (2011) Submitted

**Reinforcement learning** 

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012

#### Maximal loss

1 % choices predicted by model 5.0 5.0 0



Huys et al. (2011) Submitted

Reinforcement learning

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012





Huys et al. (2011) Submitted

Reinforcement learning

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012



Huys et al. (2011) Submitted

Reinforcement learning

Advanced Course in Computational Neuroscience, Bedlewo, Poland August 15-16 2012



# Pruning parameters

- Given the model, can now look at parameters
  - Ensure they are meaningful



# Pruning parameters

- Given the model, can now look at parameters
  - Ensure they are meaningful



# Pruning parameters

- Given the model, can now look at parameters
  - Ensure they are meaningful



#### Pruning in the brain: PAG and sgPFC



#### **Two Pavlovian influences**

# (i.e. state values directly linked to actions thorughout learning)



# Recap

- Multiple decision systems
- Multiple values
- Multiple action mechanisms
- Interactions
  - Override
  - Uncertainty
- Complex problem
- Identification via critical features