

Genela Morris Dept. of Neurobiology Haifa University gmorris@sci.haifa.ac.il

#### the role of dopamine in planning and action

#### ON NEURAL CORRELATES OF REINFORCEMENT LEARNING

#### Reinforcement learning: finding correct action by trial and error



Reinforcement learning the basics

Supervised learning – all knowing teacher, detailed feedback Reinforcement learning – scalar (correct/incorrect) feedback Unsupervised learning – self organization

# Reinforcement learning: The law of effect

"The Law of Effect is that: Of several responses made to the same situation, those which are accompanied or closely followed by satisfaction to the animal will, other things being equal, be more firmly connected with the situation, so that, when it recurs, they will be more likely to recur"

Edward Lee Thorndike (1911)

### Early attempts at modeling

- By associative rules
- Classical conditioning



## Classical conditioning

#### **The Elements:**

- **US**: Unconditioned stimulus
- **UR**: Unconditioned response
- **NS**: Neutral stimulus
- **CS**: Conditioned stimulus
  - **CS1**: Conditioned stimulus 1
  - **CS2**: Conditioned stimulus 2
- **CR**: Conditioned response



#### Properties of classical conditioning

#### (Pavlov 1927)

- Acquisition.
- Partial Reinforcement (probabilistic).
- Generalization.
- Interstimulus Interval (ISI) effects.
- Intertrial Interval (ITI) effects.

### So far...

 A simple association (coincidence, Hebbian) model can explain the phenomenon.
 Acquisition.



- Partial Reinforcemer
   (probabilistic).
- Generalization.
- Interstimulus Interval (ISI) effects.
- Intertrial Interval (ITI) effects.

#### Properties of classical conditioning

(Cnt'd)

- Conditioned Inhibition
- Latent Inhibition
- Relative validity (Wagner 1968).
- Blocking (Kamin 1968)

•

#### **CS must RELIABLY predict US**

#### Which simple association can't explain

Learning occurs not because two events co-occur, but because that co-occurrence is otherwise UNPREDICTED

## Rescorla-Wagner rule (1972)

Learning to predict reward R given stimulus U=1 Goal: Form a prediction V of the reward of the form: Where:

V=ωU

And learn to change  $\boldsymbol{\omega}$  :

 $\Delta \omega = \epsilon (R-V)U$ 

U=CS availability (0,1); V=reward prediction: R=reward availability (0,1) :  $\omega$  = weight of the connection between U and V  $\varepsilon$  = learning rate R-V = prediction error

After learning of consistent pairing:  $\omega$ =R

### Blocking with Rescorla Wagner

- Given U1, U2 and R, after U1 has been learnt:
- ω1=R
  V=ω1U1+ω2U2

R

Prediction error: R-V=0
 And no learning occurs for ω2

The real challenge faced by the nervous system is to select a motor response (to make a decision) that maximizes evolutionary fitness under conditions of uncertainty

*Glimcher, 2004* 

## Matching law (Herrnstein 1961)



#### **Experimental data**

In a VI-VI operant reward schedule The relative response rate *R* matches The relative reward rate *Rf* on that choice

$$\frac{R_1}{R_1 + R_2} = \frac{Rf_1}{Rf_1 + Rf_2}$$



### Critical problems, for control

#### 1. Exploration/exploitation





### Solutions, for control

- 1. Variability in response policy
  - Greedy ← → Random (gambling)
  - 2. Based on expected return







Weizmann systems

#### Monkeys' decisions: probability matching



#### ... whether optimal or not

• Actions are related to their consequences



# TD learning - solution for temporal credit assignment

- 1. Estimate value of current state  $(V_t = r_t + \gamma' r_{t+1} + \cdots)$ : (discounted) sum of expected rewards
- 2. Measure 'truer' value of current state: reward at present state + estimated value of next state  $(r_t + \gamma V_{t+1})$
- **3.** TD error  $\delta_t = r_t + \gamma V_{t+1} V_t$

**4.** Use TD error to improve 1  $(V_t^{k+1}=V_t^k+\eta \ \delta_t)$ where:  $V_{t=value}$  of the state reached at time t in iteration k  $r_t$  = reward given at time t;  $\eta$  = learning rate,  $\delta$  = prediction error

**TD error:**  $\delta_t = r_t + \gamma V_{t+1} - V_t$ 



#### **TD error:** $\delta_t = \gamma V_{t+1} - V_t + r_t$



#### Basal ganglia - anatomy



#### Intracranial self stimulation



#### The midbrain dopamine system



# Dopamine and acetylcholine meet in the striatum



#### Facts to remember (1)

- Basal ganglia receive cortical input
- Basal ganglia project to frontal cortex
- Dopamine and acetylcholine localization

#### The midbrain dopamine system



## Probabilistic instrumental conditioning task





### Dopamine population response-



vveizmann systems

time (ms)

#### **Dopamine population response**reward



# Dopamine population response – reward omission



#### Instrumental conditioning - results

- Responses to visual cue are correlated with future reward probability
- Responses to reward are inversely correlated with reward probability
- Responses to reward omission are indifferent to reward probability
- Dopamine neurons provide an accurate TD signal (but only in the positive domain)

# ... and it can cause long term plasticity of cortico-striatal synapses



# ... and it can cause long term plasticity of cortico-striatal synapses



Weizmann system Shen et al., 2008

#### Facts to remember 2

- DA neurons provide a TD error signal
- To the cortico (state) striatal (action) synapses
- And DA modulates synaptic plasticity

### **Control - Adding action**



#### The agent has to:

- Learn to predict reinforcement
- Know the state-action-state transitions
   *policy*

state value behavioural

#### Alternatively – reinforce action







No explicit knowledge about the future

#### Solution 1: actor/critic networks



# How can the dopamine signal contribute to decision behaviour?

Long term policy-shaping effect

through synaptic plasticity



Immediate effect on action















#### Lost in translation?



### Monkeys' decisions: shaping by dopamine R - DA response rate (spikes/s) 25 50 75 reward probability - right 0.5R<sup>2</sup>= 0.930 0 0.5 0 Dright/(Dright + Dleft)

#### Dopamine neurons during decision



#### Are DA neurons aware of future choice

Hi (explore)Lo (exploit)



#### The learning is of state-action values



### Adding an internal model

Ctx D1/5 STR





- No explicit knowledge about the future is necessary
- But what about planning?

# Reinforcement devaluation – evidence for model



Nature Reviews | Neuroscience

## Planning and decision making



#### Model based learning

- Creating an internal model of the environment
- Supported by PFC and Nac
- Where is the model?
- How is it learnt?

#### A cognitive map of space in the brain?



Edward Tolman



- Not everything is stimulus response
- Rather, there are internal representations
- Cognitive maps as an example of an internal/mental representation

Tolman EC (1948) Cognitive maps in rats and men. Psychol Rev 55:189-208.

### **Model based learning**

We believe that in the course of learning something like a field map of the environment gets established in the rat's brain... Although we admit that the rat is bombarded by stimuli, we hold that his nervous system is surprisingly selective as to which of these stimuli it will let in at any given time... Rather, the incoming impulses are usually worked over and elaborated in the central control room into a tentative, cognitive-like map of the environment. And it is this tentative map, indicating routes and paths and environmental relationships, which finally determines what responses, if any, the animal will finally release.

> Cognitive maps of rats and men (1948) The Psychological Review, 55(4), 189-208



Edward C. Tolman

### **Model based learning**

We believe that in the course of learning something like a field map of the environment gets established in the rat's brain... Although we admit that the rat is bombarded by stimuli, we hold that his nervous system is surprisingly selective as to which of these stimuli it will let in at any given time... Rather, the incoming impulses are usually worked over and elaborated in the central control room into a tentative, cognitive-like map of the environment. And it is this tentative map, indicating routes and paths and environmental relationships, which finally determines what responses, if any, the animal will finally release.

> Cognitive maps of rats and men (1948) The Psychological Review, 55(4), 189-208



Edward C. Tolman

### **Model based learning**

We believe that in the course of learning something like a field map of the environment gets established in the rat's brain... Although we admit that the rat is bombarded by stimuli, we hold that his nervous system is surprisingly selective as to which of these stimuli it will let in at any given time... Rather, the incoming impulses are usually worked over and elaborated in the central control room into a tentative, cognitive-like map of the environment. And it is this tentative map, indicating routes and paths and environmental relationships, which finally determines what responses, if any, the animal will finally release.

> Cognitive maps of rats and men (1948) The Psychological Review, 55(4), 189-208



Edward C. Tolman