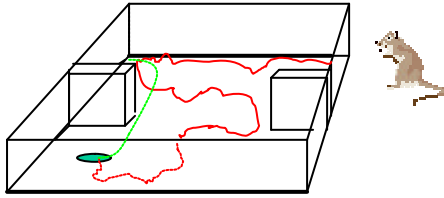
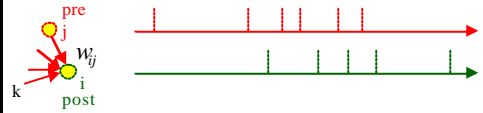


**Unsupervised vs. reinforcement learning
(via a model of rat navigation)**



LeN

Hebbian Learning

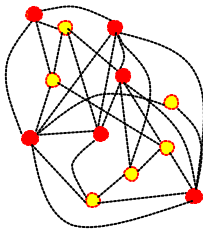


When an axon of cell j repeatedly or persistently takes part in firing cell i , then j 's efficiency as one of the cells firing i is increased

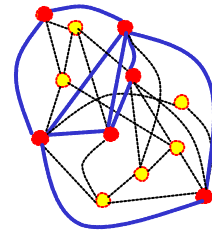
Hebb, 1949

- local rule
- simultaneously active (correlations)

Hebbian Learning



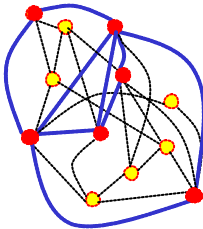
Hebbian Learning



item memorized

Hebbian Learning

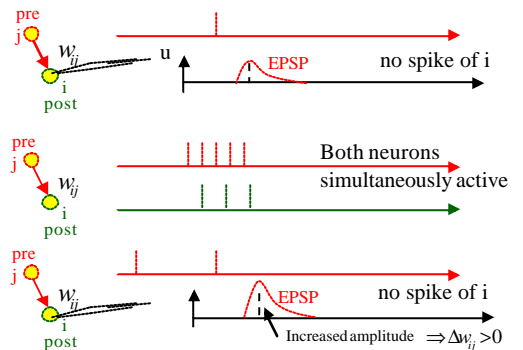
Recall:
Partial info



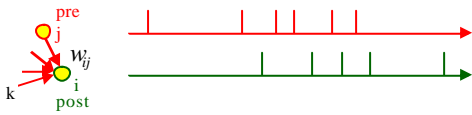
item recalled

→ Associative memory

Hebbian Learning in experiments (schematic)



Hebbian Learning



When an axon of cell j repeatedly or persistently takes part in firing cell i , then j 's efficiency as one of the cells firing i is increased

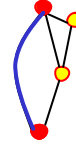
Hebb, 1949

- local rule
- simultaneously active (correlations)

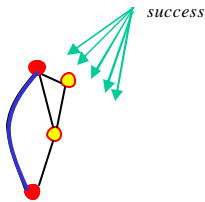
Rate model:

active = high rate = many spikes per second

Hebbian Learning = unsupervised learning



Reinforcement Learning = reward + Hebb

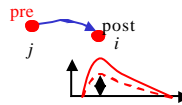


Classification of plasticity: unsupervised vs reinforcement

LTP/LTD/Hebb

Theoretical concept

- passive changes
- exploit statistical correlations



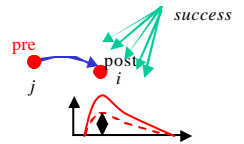
Functionality

- useful for development
(wiring for receptive fields)

Reinforcement Learning

Theoretical concept

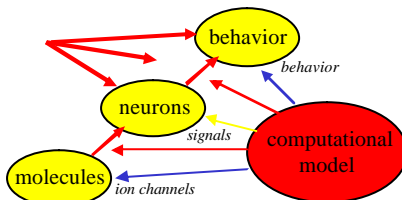
- conditioned changes
- maximise reward



Functionality

- useful for learning
a new behavior

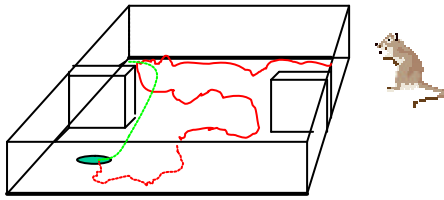
Introduction to reinforcement learning (via a model of rat navigation)



Introduction to reinforcement learning (via a model of rat navigation)

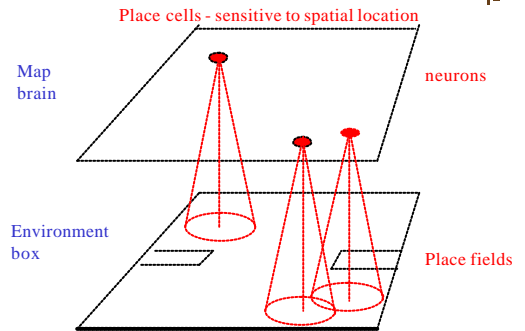
- -Basics of rat navigation
- Place cells and Rat hippocampus
- A model of spatial representation
- Learning to find the goal location
- Reward based learning
- Reinforcement learning theory
- Eligibility traces
- the full model: behavioral experiments

Biological Principles of Learning: spatial learning

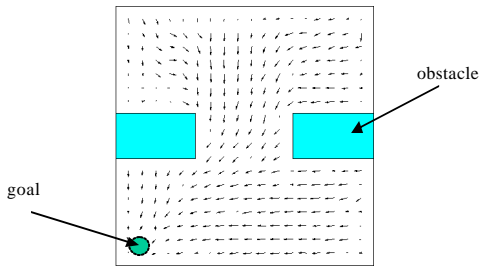


LeN

Spatial representation



Use map



LeN

SPATIAL REPRESENTATION

- Model of place cells

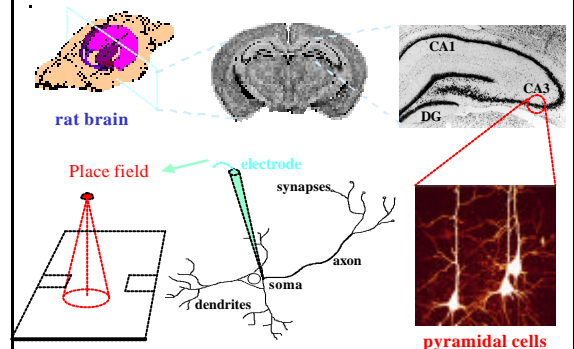
GOAL LEARNING

- Reward-based learning system

Introduction to reinforcement learning (via a model of rat navigation)

- Basics of rat navigation
- -Place cells and Rat hippocampus
- A model of spatial representation
- Learning to find the goal location
- Reward based learning
- Reinforcement learning theory
- Eligibility traces
- the full model: behavioral experiments

Neurophysiology of the Rat Hippocampus



pyramidal cells

Hippocampal Place Cells

Depends on

- visual cues
- works also in the dark

place field
O'Keefe
Place field

The Neural Model: Place Cells

Hippocampus CA3-CA1

vision place cells
Visual Processing
Visual Stimuli

integration place cells
Path Integrator
Internal Stimuli

The Neural Model: Place Cells

Hippocampus

recording

place fields

The Neural Model: Place Cells

CA3-CA1 population firing

place field center

Center of mass of place cell activity

✓ **SPATIAL REPRESENTATION**

- Model of place cells

→ **GOAL LEARNING**

- Reward-based learning system

Reward-based Action Learning

Hippocampal place cells = fuzzy discretisation of continuous space

Action Learning
NA

Spatial representation
Place Cells

External stimuli Internal stimuli

Assign 'value' to states and actions (Bellman equation/dyn. Programming)

reinforcement-learning in continuous space (Q-learning)

goal

Reward-based Action Learning

Connection reinforced if action a at state s successful

Success = reward - exp reward

$$\Delta w_{ij} = h d_t e_{ij}$$

$$e_t(s) = r(s) + g \cdot e_{t-1}(s)$$

State = activity $r(s)$

Action a = north

Reward-based Action Learning

Connection reinforced if action a at state s successful

Success = reward - exp reward

$$\Delta w_{ij} = h d_t e_{ij}$$

Dopamine neurons (W. Schultz)

Reward-based Action Learning

Connection reinforced if action a at state s successful

Success signal

Local rule, conditioned on global success

Introduction to reinforcement learning (via a model of rat navigation)

- Basics of rat navigation
- Place cells and Rat hippocampus
- A model of spatial representation
- Learning to find the goal location
- Reward based learning (basic ideas)
- Reinforcement learning theory
- Eligibility traces
- the full model: behavioral experiments

Reward-based Action Learning

Connection reinforced if action a at state s successful

$Q(s,a) = \text{expected reward}$

$$Q(s,a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a$$

Iterative Update

$$DQ(s,a) = h [r - Q(a)]$$

Blackboard:

- $s = \text{state}$
- $a = \text{action}$
- $R_{s \rightarrow s'}^a = \text{reward}$
- $s' = \text{new action}$

Exercise now: Iterative update

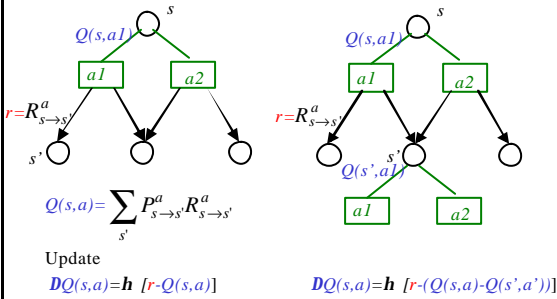
- Show that an empirical evaluation of $Q(s,a)$ by averaging the rewards for action a over k or $(k+1)$ trials, leads to an iterative update rule of the form $DQ(s,a) = h [r - Q(a)]$
- Calculate η .
- Give an intuitive explanation of the update rule

$$Q(s,a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a \quad \text{expected reward}$$

Reward-based Action Learning

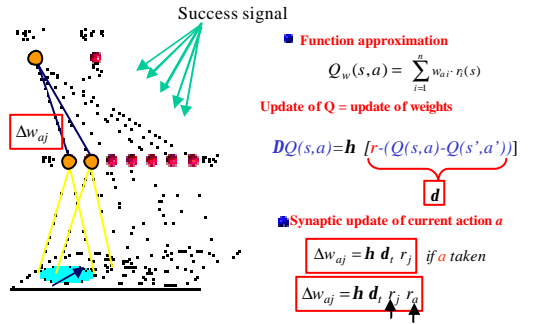
Connection reinforced if
action a at state s successful

Blackboard:



Reward-based Action Learning

Connection reinforced if
action a at state s successful

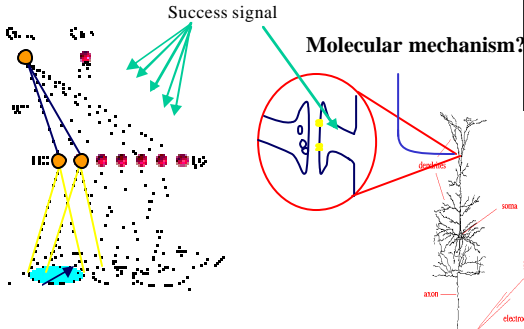


Reward-based Action Learning

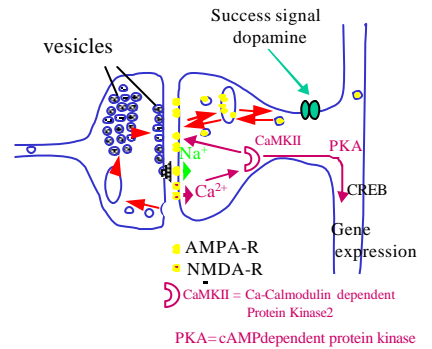
Connection reinforced if
action a at state s successful

Success signal

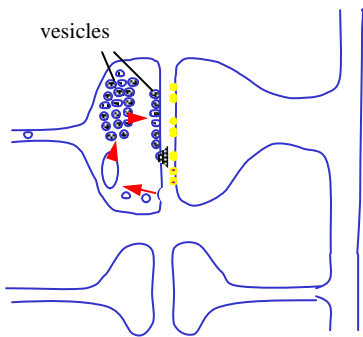
Molecular mechanism?



Changes in synaptic connections



Changes in synaptic connections



Reward-based Learning TD ()

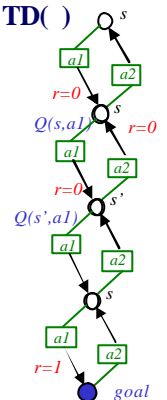
SARSA

$$DQ(s,a) = h [r - (Q(s,a) - Q(s',a'))]$$

policy for action choice:

Pick most often action
 $a_t^* = \arg \max_a Q_a(s,a)$

Blackboard:



Exercise now

Update of Q values in SARSA

$$DQ(s,a) = h [r - Q(s,a) - Q(s',a')]$$

policy for action choice:

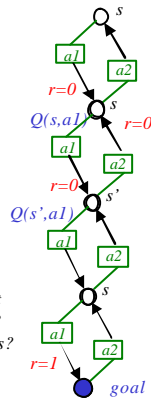
Pick most often action

$$a_t^* = \arg \max_a Q_a(s, a)$$

Consider a linear sequence of states. Reward only at goal. Actions are up or down.

- Initialise Q values at 0. Start at top. How do Q values develop?
- Q values after 3 complete trials?

goal



Problem: learning is slow

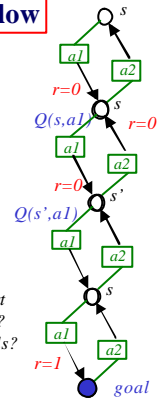
- Slow diffusion of information across several states

$$a_t^* = \arg \max_a Q_a(s, a)$$

Consider a linear sequence of states. Reward only at goal. Actions are up or down.

- Initialise Q values at 0. Start at top. How do Q values develop?
- Q values after 3 complete trials?

goal



Introduction to reinforcement learning (via a model of rat navigation)

- Basics of rat navigation
- Place cells and Rat hippocampus
- A model of spatial representation
- Learning to find the goal location
- Reward based learning (basic ideas)
- Reinforcement learning theory
- -Eligibility traces
- the full model: behavioral experiments

Reward-based Learning TD()

TD error in SARSA

$$d_t = R_{t+1} - [Q_a(s', a') - g \cdot Q_a(s, a)]$$

Function approximation

$$Q_a(s, a) = \sum_{s'} w_{s'}^a \cdot r_t(s)$$

policy for action choice:

Pick most often action

$$a_t^* = \arg \max_a Q_a(s, a)$$

Eligibility trace (memory at synapse):

$$e_{aj}(t) = r_j r_a + \begin{cases} g e_{aj}(t) & \text{if exploiting} \\ 0 & \text{if exploring} \end{cases}$$

Synaptic update

$$\Delta w_{aj} = h d_t e_{aj}$$

pre post

memory



Reward-based Action Learning

Connection reinforced if action a at state s successful

Success signal

Learning Rule

$$\Delta w_{aj} = h d_t e_{aj}$$

Spatial Representation



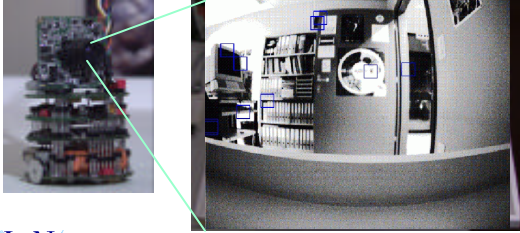
Introduction to reinforcement learning (via a model of rat navigation)

- Basics of rat navigation
- Place cells and Rat hippocampus
- A model of spatial representation
- Learning to find the goal location
- Reward based learning (basic ideas)
- Reinforcement learning theory
- Eligibility traces
- -the full model: behavioral experiments

Validating the Model

The KHEPERA mobile miniature robot

Experimental arena

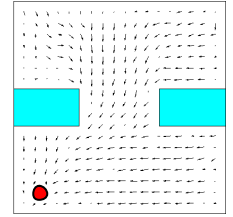
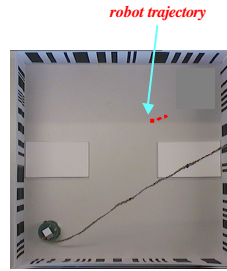


422 x 316 pixels

LeN

Open-field Navigation Experiments

(Biol. Cybern., 2000)



Navigation map after 20 training trials

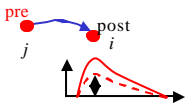
LeN

Classification of plasticity: unsupervised vs reinforcement

LTP/LTD/Hebb

Theoretical concept

- passive changes
- exploit statistical correlations



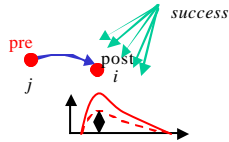
Functionality

- useful for development
- (wiring for receptive fields)

Reinforcement Learning

Theoretical concept

- conditioned changes
- maximise reward

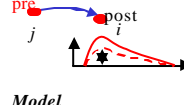


Functionality

- useful for learning a new behavior

Plasticity models: unsupervised vs reinforcement

STDP/Hebb



Model

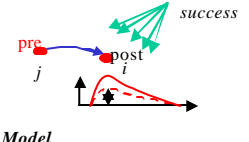
- STDP (see above)

$$\Delta w_{ij} \propto \text{pre} \cdot \text{post} + \text{other terms}$$

Reinforcement Learning

theoretical Protocol

- maximise reward

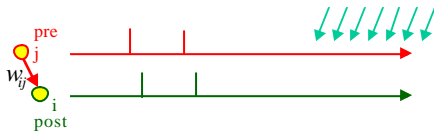


Model

- spike-based model?

$$\Delta w_{ij} \propto \text{success} \cdot (\text{pre} \cdot \text{post} + \text{other terms})$$

Timing issues in Reinforcement Learning



Success signal is delayed

Spike time scale 1-10 ms

Reward delay 100-5000 ms

→ Need memory trace (eligibility trace)