





Reinforcement Learning: From neural processes modelling to Robotics applications

Mehdi Khamassi (CNRS, ISIR-UPMC, Paris)

30 January 2015

Michèle Sebag's course @ Univ. Orsay

Team: Architectures and Models of Adaptation and Cognition (AMAC)

TDRL Model Dopamine TD Applications

slide # 2 / 141

@ The Institute of Intelligent Systems and Robotics (ISIR), UPMC-CNRS



RESEARCH INTERESTS

- TDRL Model Dopamine TD Applications Model-based R1
 - slide # 3 / 141
- **Decision-making**: choice at each moment of the most appropriate behavior for an agent's survival, to solve a task.
- **Reinforcement Learning** (by trial/error): adaptation of this choice to maximize a particular reward function.
- **Complex problems**: noise, partial representation of states, non stationarity of the environment.
- Modular/hierarchical structure of different learning levels, enabling a better flexibility and autonomy of decision in animals and robots.

Global organization of learning in the brain (according to Doya 2000)



Dopamine TD Applications

slide # 4 / 141



TDRL Model Dopamine TD Applications

slide # 5 / 141

1. Model-free Reinforcement Learning

- Temporal-Difference RL Algorithm
- Dopamine activity
- Wide application to Neuroscience of decision-making

2. Model-based Reinforcement Learning

- Off-line learning / Replay during sleep
- Dual-system RL
- Online parameters tuning (meta-learning)
- Link with Neurobehavioral data
- Applications to Robotics

TDRL Model

TD Applications

slide # 6 / 141

REINFORCEMENT LEARNING & DOPAMINE ACTIVITY

THE ACTOR-CRITIC MODEL

Sutton & Barto (1998) Reinforcement Learning: An Introduction

TDRL Model

Dopamine TD Applications

<u>slide # 7 / 141</u>



The Actor learns to select actions that maximize reward.

The Critic learns to predict reward (its value V).

A reward prediction error constitutes the reinforcement signal.

Dopamine TD Applications

slide # 8 / 141

. Learning from delayed reward





TDRL Model

Dopainine D Applications

slide # 9 / 141

. Learning from delayed reward



slide # 10 / 141



slide # 11 / 141

Temporal-Difference (TD) learning

•



TDRL Model

Dopamine FD Applications

slide # 12 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

$$V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$$

learning rate (=0.9)

TDRL Model

TD Applications

slide # 13 / 141



0 = 0 + 0 - 0 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (

discount factor (=0.9)

 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

0 = 0 + 0.9 * 0

TDRL Model

Dopamine TD Applications

slide # 14 / 141



1 = 1 + 0 - 0 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (=0.9)

0.9 = 0 + 0.9 * 1 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

TD Applications

slide # 15 / 141



1 = 1 + 0 - 0 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (=0.9)

0.9 = 0 + 0.9 * 1 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

TDRL Model

TD Applications

slide # 16 / 141



$$0 = 0 + 0 - 0$$

$$\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$$

discount factor

discount factor (=0.9)

 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

0 = 0 + 0.9 * 0

TDRL Model

Dopamine FD Applications

slide # 17 / 141



0.81 = 0 + 0.9 * 0.9 - 0 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (=0.9) 0.72 = 0 + 0.9 * 0.81 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

Dopamine TD Applications

slide # 18 / 141



0.81 = 0 + 0.9 * 0.9 - 0 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (=0.9)

0.72 = 0 + 0.9 * 0.81 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

TDRL Model

TD Applications

slide # 19 / 141



0.1 = 1 + 0 - 0.9 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (=0.9)

0.99 = 0.9 + 0.9 * 0.1 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

TD Applications

slide # 20 / 141



0.1 = 1 + 0 - 0.9 $\delta_{t+1} = r_{t+1} + \gamma \cdot V(s_{t+1}) - V(s_t)$ discount factor (=0.9)

0.99 = 0.9 + 0.9 * 0.1 $V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$ learning rate (=0.9)

TDRL Model

Dopamine TD Applications

slide # 21 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

 $V(s_{t}) = V(s_{t}) + \alpha \cdot \delta_{t+1}$ learning rate (=0.1) usually small for stability

TDRL Model

Dopamine TD Applications

slide # 22 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

TDRL Model

Dopamine TD Applications

slide # 23 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

$$V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$$

learning rate (=0.1)

TDRL Model

Dopamine TD Applications

slide # 24 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

TDRL Model

Dopamine TD Applications

slide # 25 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

TDRL Model

Dopamine FD Applications

slide # 26 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

TDRL Model Dopamine TD Applications

slide # 27 / 141

How can the agent learn a policy?

How to learn to perform the right actions

TDRL Model Dopamine TD Applications

slide # 28 / 141

How can the agent learn a policy?

How to learn to perform the right actions

- S: state space
- A : action space
- Policy function $\pi: S \longrightarrow A$

What we have learned so far:

Value function $V: S \longrightarrow \mathbb{R}$

slide # 29 / 141

How can the agent learn a policy?

How to learn to perform the right actions

a solution: parallely update a policy and a value function



$$\mathbf{V}(\mathbf{s}_t) = \mathbf{V}(\mathbf{s}_t) + \boldsymbol{\alpha} \cdot \boldsymbol{\delta}_{t+1}$$

slide # 30 / 141

How can the agent learn a policy?

How to learn to perform the right actions

other solution: learning Q-values (qualities)

 $Q: (S,A) \longrightarrow \mathbb{R}$

Q-table:

state / action	a1 : North	a2 : South	a3 : East	a4 : West
s1	0.92	0.10	0.35	0.05
s2	0.25	0.52	0.43	0.37
s3	0.78	0.9	1.0	0.81
s4	0.0	1.0	0.9	0.9

slide # 31 / 141

How can the agent learn a policy?

How to learn to perform the right actions

other solution: learning Q-values (qualities)

 $Q: (S,A) \longrightarrow \mathbb{R}$ Q-table:



state / action	a1 : North	a2 : South	a3 : East	a4 : West
s1	0.92	0.10	0.35	0.05
s2	0.25	0.52	0.43	0.37
s3	0.78	0.9	1.0	0.81
s4	0.0	1.0	0.9	0.9

slide # 32 / 141

How can the agent learn a policy?

How to learn to perform the right actions

other solution: learning Q-values (qualities)

 $Q: (S,A) \longrightarrow \mathbb{R}$ Q-table:

state / action	a1 : North	a2 : South	a3 : East	a4 : West
s1	0.92	0.10	0.35	0.05
s2	0.25	0.52	0.43	0.37
s3	0.78	0.9	1.0	0.81
s4	0.0	1.0	0.9	0.9

$$P(a) = \frac{\exp(\beta \cdot Q(s,a))}{\sum_{b} \exp(\beta \cdot Q(s,b))}$$

The β parameter regulates the exploration – exploitation trade-off.

TDRL Model Dopamine

slide # 33 / 141

ACTOR-CRITIC

 $V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$

State-dependent Reward Prediction Error

(independent from the action)

slide # 34 / 141

ACTOR-CRITIC

 $V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$

State-dependent Reward Prediction Error

(independent from the action)



TD Applications

ACTOR-CRITIC

 $V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$

SARSA

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

Reward Prediction Error dependent on the action

chosen to be performed next

slide # 36 / 141

ACTOR-CRITIC

 $V(s_t) \leftarrow V(s_t) + \alpha[r_{t+1} + \gamma V(s_{t+1}) - V(s_t)]$

SARSA

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

• Q-LEARNING

 $Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[r_{t+1} + \gamma \max_{a \in A} Q(s_{t+1}, a) - Q(s_t, a_t)]$

Reward Prediction Error dependent on the best action
TDRL Model Dopamine D Applications

slide # 37 / 141

Links with biology

Activity of dopaminergic neurons

TD-learning explains classical conditioning (predictive learning)



The unconditioned stimulus is repeatedly presented just after the neutral stimulus. The unconditioned stimulus continues to produce an unconditioned response. The neutral stimulus alone now produces a conditioned response (CR), thereby becoming a conditioned stimulus (CS).

Taken from Bernard Balleine's lecture at Okinawa Computational Neuroscience Course (2005).

TDRL Model Dopamine

slide # 39 / 141





TDRL Model Dopamine TD Applications

slide # 40 / 141



TDRL Model Dopamine TD Applications

slide # 41 / 141



Dopamine TD Applications

slide # 42 / 141



The Actor-Critic model

Dopamine

slide # 43 / 141



The Actor-Critic model

TDRL Model Dopamine TD Applications

slide # 44 / 141



Montague et al. (1996); Suri & Schultz (2001) Daw (2003); Bertin et al. (2007).

The Actor-Critic model

Dopamine

slide # 45 / 141



TDRL Model Dopamine D Applications

slide # 46 / 141

Wide application of RL models to model-based analyses of behavioral and physiological data during decision-making tasks

Typical probabilistic decision-making task

<u>slide # 47 / 141</u>

Dopamine



Niv et al. (2006), commentary about the results presented in Morris et al. (2006) Nat Neurosci.

Katie Ris

Typical probabilistic decision-making task

slide # 48 / 141

Dopamine



Niv et al. (2006), commentary about the results presented in Morris et al. (2006) Nat Neurosci.

Typical probabilistic decision-making task

slide # 49 / 141

Dopamine



Niv et al. (2006), commentary about the results presented in Morris et al. (2006) Nat Neurosci.

Model-based analysis of brain data

TDRL Model Dopamine TD Applications

slide # 50 / 141

Sequence of observed trials : Left (Reward); Left (Nothing); Right (Nothing); Left (Reward); ...



cf. travail de Mathias Pessiglione (ICM)

ou Giorgio Coricelli (ENS)

Model-based analysis Work by Jean Bellot (PhD student)



Model-based analysis My post-doc work

TDRL Model Dopamine TD Applications

slide # 52 / 141



- Analysis of single neurons recorded in the monkey dorsolateral prefrontal cortex and anterior cingulate cortex
- Correlates of prediction errors? Action values? Level of control/exploration?

Khamassi et al. (2013) Prog Brain Res; Khamassi et al. (in revision)

Model-based analysis My post-doc work

TDRL Model Dopamine TD Applications

slide # 53 / 141



Multiple regression analysis with bootstrap

Khamassi et al. (2013) Prog Brain Res; Khamassi et al. (in revision)

TDRL Model Dopamine TD Applications

slide # 54 / 141

This works well, but...

Most experiments are single-step

All these cases are discrete

•

•

•

- Very small number of states, actions
- We supposed a perfect state identification

TDRL Model Dopamine TD Applications

slide # 55 / 141



TD-Learning model applied to spatial navigation behavior learning in a robot performing the bio-inspired plus-maze task

Khamassi et al. (2005). Adaptive Behavior. Khamassi et al. (2006). Lecture Notes in Computer Science

TDRL Model Dopamine TD Applications

slide # 56 / 141

Coordination by a self-organizing map



TDRL Model Dopamine TD Applications

slide # 57 / 141





Hand-tuned

Autonomous

Random

TDRL Model Dopamine TD Applications

slide # 58 / 141

Two methods :



Autonomous

1. Self-Organizing Maps (SOMs)

2. specialization based on performance (tests modules' capacity for state prediction) Baldassarre (2002); Doya et al. (2002). Within a particular subpart of the maze, only the module with the most accurate reward prediction is trained. Each module thus becomes an expert responsible for learning in a given task subset.

TDRL Model Dopamine TD Applications

slide # 59 / 141



TDRL Model Dopamine TD Applications

slide # 60 / 141

Nb of iterations required

(Average performance during the second half of the experiment)

1. hand-tuned	94
2. specialization based on performance	3,500
3. autonomous categorization (SOM)	404
4. random robot	30,000





TDRL Model Dopamine TD Applications



Nb of iterations required

(Average performance during the second half of the experiment)





TDRL Model Dopamine TD Applications

slide # 62 / 141

1. Model-free Reinforcement Learning

- Temporal-Difference RL Algorithm
- Dopamine activity
- Wide application to Neuroscience of decision-making

2. Model-based Reinforcement Learning

- Off-line learning / Replay during sleep
- Dual-system RL
- Online parameters tuning (meta-learning)
- Link with Neurobehavioral data
- Applications to Robotics

TDRL Model Dopamine TD Applications Model-based RL slide # 63 / 141

Off-learning (Model-based RL) and hippocampal & prefrontal cortex activity replay during sleep

TDRL Model Dopamine TD Applications Model-based RL slide # 64 / 141



$$\delta_{t+1} = \mathbf{r}_{t+1} + \gamma \cdot \mathbf{V}(\mathbf{s}_{t+1}) - \mathbf{V}(\mathbf{s}_t)$$

discount factor (=0.9)

$$V(s_t) = V(s_t) + \alpha \cdot \delta_{t+1}$$

learning rate (=0.1)

TRAINING DURING SLEEP

TDRL Model Dopamine TD Applications Model-based RL

slide # 65 / 141







Method in Artificial Intelligence: Off-line Dyna-Q-learning (Sutton & Barto, 1998) To incrementally learn a model of transition and reward functions, then plan within this model by updates "in the head of the agent" (Sutton, 1990).



TDRL Model Dopamine TD Applications Model-based RL slide # 67 / 141

s : state of the agent (\bullet)



Dopamine TD Applications Model-based RL slide # 68 / 141

s : state of the agent (\bullet)



TD Applications Model-based RL slide # 69 / 141

- s : state of the agent (\bullet)
- a : action of the agent (go east)



TD Applications Model-based RL slide # 70 / 141

- s : state of the agent (\bullet)
- a : action of the agent (go east)
- stored transition function T: $proba(\longrightarrow) = 0.9$ $proba(\swarrow) = 0.1$ $proba(\searrow) = 0$



Model-based RL s : state of the agent (\bullet) a : action of the agent (go east) maxQ=0.3maxQ=0.9maxQ=0.7 stored transition function T: $proba(\longrightarrow) = 0.9$ proba(>) = 0.1 $proba(\searrow) = 0$ $\mathcal{Q}(s,a) \leftarrow \mathcal{R}(s,a) + \gamma \sum \mathcal{T}(s'|s,a) \max \mathcal{Q}(s',a')$ s' $0.9*0.7 + 0.1*0.9 + 0*0.3 + \dots$ 0.6

Model-based Reinforcement Learning

TDRL Model Dopamine TD Applications Model-based RL slide # 72 / 141

No reward prediction error!

Only:

- Estimated Q-values
- Transition function
- Reward function

This process is called Value Iteration or Dynamic prog.

$$\mathcal{Q}(s, a) \leftarrow \mathcal{R}(s, a) + \gamma \sum_{s'} \mathcal{T}(s'|s, a) \max_{a'} \mathcal{Q}(s', a')$$
TDRL Model Dopamine TD Applications Model-based RL slide # 73 / 141

Links with neurobiological data

Activity of hippocampal place neurons

Hippocampal place cells

TDRL Model Dopamine TD Applications Model-based RL slide # 74 / 141

a Cell 1 11.00.0001 Cell 2 U. R. HILL HIL Cell 3 1111001001 Cell 4 || || || || || || b Nature Reviews | Neuroscience

Nakazawa, McHugh, Wilson, Tonegawa (2004) Nature Reviews Neuroscience

Hippocampal place cells

Model-based RL slide # 75 / 141

PRE RUN POST

· Reactivation of hippocampal place cells during sleep (Wilson & McNaughton, 1994, Science)



Hippocampal place cells

•

TDRL Model Dopamine TD Applications Model-based RL slide # 76 / 141



Forward replay of hippocampal place cells during sleep (sequence is compressed 7 times) (Euston et al., 2007, Science)

Sharp-Wave Ripple (SWR) events

TDRL Model Dopamine TD Applications Model-based RL slide # 77 / 141

 "Ripple" events = irregular bursts of population activity that give rise to brief but intense highfrequency (100-250 Hz) oscillations in the CA1 pyramidal cell layer.





Selective suppression of SWRs impairs spatial memory

TDRL Model Dopamine TD Applications Model-based RL

slide # 78 / 141



 Girardeau G, Benchenane K, Wiener SI, Buzsáki G, Zugaro MB (2009) Nat Neurosci.

Contribution to decision making (forward planning) and evaluation of transitions



Figure 1. The multiple-T maze. The task consists of four T choice points with food reward available at two sites on each return rail. Only feeders on one side of the track were rewarded in each session.



Johnson & Redish (2007) J Neurosci

Model-based RL

slide # 79 / 141

SUMMARY OF NEUROSCIENCE DATA

Replay their sequential activity during sleep (Foster & Wilson, 2006; Euston et al., 2007; Gupta et al., 2010)

- Performance is impaired if this replay is disrupted (Girardeau, Benchenane et al. 2012; Jadhav et al. 2012)
- Only task-related replay in PFC (Peyrache et al., 2009)
- Hippocampus may contribute to model-based navigation strategies, striatum to model-free navigation strategies (Khamassi & Humphries, 2012)

Applications to robot off-line learning Work of Jean-Baptiste Mouret et al. @ ISIR TDRL Model Dopamine TD Applications Model-based RL slide # 81 / 141

How to recover from damage without needing to identify the damage?



Applications to robot off-line learning Work of Jean-Baptiste Mouret et al. @ ISIR TDRL Model Dopamine TD Applications Model-based RL slide # 82 / 141

The reality gap

Self-model vs reality: how to use a simulator?



Solution: Learn a transferability function (how well does the simulation match reality?) with SVM or neural networks.

Idea: the damage is a large reality gap.

Koos, Mouret & Doncieux. IEEE Trans Evolutionary Comput 2012

Applications to robot off-line learning Work of Jean-Baptiste Mouret et al. @ ISIR TDRL Model Dopamine TD Applications Model-based RL slide # 83 / 141

Experiments



Koos, Cully & Mouret. Int J Robot Res 2013

TDRL Model Dopamine TD Applications MMeta-Learning

slide # 84 / 141

META-LEARNING (regulation of decision-making) 1. Dual-system RL coordination 2. Online parameters tuning

Multiple decision systems

TDRL Model Dopamine TD Applications Model-based RL slide # 85 / 141

Model-based system Skinner box (instrumental conditioning) Model-free sys. b а S₀ S₀ Initial state Initial state Press Enter Press Enter lever magazine lever magazine Q = 1Q = 0s, Food delivered S, S2 Press Enter lever magazine Food delivered No reward Q = 0Q = 1Enter Press R = 0magazine lever S2 No reward Q = 0s, S3 S₃ Food obtained Food obtained No reward R = 0R = 1Q = 1

(Daw Niv Dayan 2005, Nat Neurosci)

Behavior is initially model-based and becomes modelfree (habitual) with overtraining.



Dopamine TD Applications Model-based RL

<u>slide # 86 / 141</u>

Yin et al. 2004; Balleine 2005; Yin & Knowlton 2006

TDRL Model Dopamine TD Applications Model-based RL

slide # 87 / 141



Yin et al. 2004; Balleine 2005; Yin & Knowlton 2006

TDRL Model Dopamine TD Applications Model-based RL

slide # 88 / 141





TDRL Model Dopamine TD Applications Model based BL

Model-based RL slide # 89 / 141

Yin et al. 2004; Balleine 2005; Yin & Knowlton 2006

Model-free vs model-based: outcome sensitivity

Switch with

experience

computational

[reduce

load]

TDRL Model Dopamine TD Applications Model-based RL

slide # 90 / 141



Change R: slow to update

Habitual



Change R: fast to update

Goal-directed

Daw et al 2005 Nat Neurosci

TDRL Model Dopamine TD Applications Model-based RL slide # 91 / 141

Keramati et al. (2011): extension of the Daw 2005 model with a speed-accuracy trade-off arbitration criterion.



Progressive shift from model-based navigation to model-free navigation

TDRL Model Dopamine TD Applications Model-based RL slide # 92 / 141







Model-based and model-free navigation strategies

TDRL Model Dopamine TD Applications Model-based RL slide # 93 / 141

Model-based navigation



Model-free navigation



Benoît Girard 2010 UPMC lecture

Old behavioral evidence for Place-based model-based RL

TDRL Model Dopamine TD Applications Model-based RL slide # 94 / 141



Martinet et al. (2011) model applied to the Tolman maze

Old behavioral evidence for Place-based model-based RL

TDRL Model Dopamine TD Applications Model-based RL slide # 95 / 141

Martinet et al. (2011) model applied to the Tolman maze

MULTIPLE NAVIGATION STRATEGIES IN THE RAT

TDRL Model Dopamine TD Applications

Model-based RL slide # 96 / 141





Devan and White, 1999

MULTIPLE DECISION SYSTEMS IN A NAVIGATION MODEL

Model-based

system

(hippocampal

place cells)



Model-free system (basal ganglia)

Work by Laurent Dollé:

Dollé et al., 2008, 2010, submitted

MULTIPLE NAVIGATION STRATEGIES IN A TD-LEARNING MODEL

TDRL Model Dopamine TD Applications Model-based RL slide # 98 / 141

Task with a cued platform (visible flag) changing location every 4 trials



Task of Pearce et al., 1998





Model:

Dollé et al., 2010

PSIKHARPAX ROBOT



Work by: Ken Caluwaerts (2010) Steve N'Guyen (2010) Mariacarla Staffa (2011) Antoine Favre-Félix (2011)



Caluwaerts et al. (2012) Biomimetics & Bioinspiration

PSIKHARPAX ROBOT

Model-based RL slide # 100 / 141

Planning strategy only

Planning strategy + Taxon strategy



(a)

Caluwaerts et al. (2012) Biomimetics & Bioinspiration

CURRENT APPLICATIONS TO THE PR2 ROBOT





Erwan Renaudo

Omar Islas Ramirez





CURRENT APPLICATIONS TO HUMAN-ROBOT INTERACTION

TDRL Model Dopamine TD Applications Model-based RL

slide # 102 / 141

Travaux de : Erwan Renaudo Collaboration : Alami et al (LAAS)



(a) Initial state

(b) 3d model view of initial state

Task: Clean the table
Current state: A priori given action plan (right image)
Goal: Autonomous learning by the robot



Sign-trackers



Goal-trackers







Dopamine TD Applications MMeta-Learning slide # 104 / 141

Sign-trackers Goal-trackers





Fast Scan Cyclic Voltammetry (FSCV) in the ventral striatum.

TDRL Model Dopamine TD Applications MMeta-Learning slide # 105 / 141

Sign-trackers

Goal-trackers



Fast Scan Cyclic Voltammetry (FSCV) in the ventral striatum.

TDRL Model Dopamine TD Applications MMeta-Learning slide # 106 / 141



Systemic injection of flupentixol prior to each session.

Dopamine Dopamine D Applications

MMeta-Learning slide # 107 / 141

Computational model



Lesaint, Sigaud, Flagel, Robinson, Khamassi (2014) PLOS Computational Biology.

Dopamine Dopamine

Computational model

MMeta-Learning slide # 108 / 141



Lesaint, Sigaud, Flagel, Robinson, Khamassi (2014) PLOS Computational Biology.
Modelling the task as a Markov Decision Process



TDRL Model Dopamine D Applications

MMeta-Learning slide # 110 / 141



$$\begin{cases} \mathcal{P}_{goL}(s_1) &= (1-\omega) \quad (\mathcal{Q}_{goL}(s_1) - \max_{a'} \mathcal{Q}_{a'}(s_1)) &+ \omega \quad \mathcal{V}(L) \\ \\ \mathcal{P}_{goM}(s_1) &= (1-\omega) \quad (\mathcal{Q}_{goM}(s_1) - \max_{a'} \mathcal{Q}_{a'}(s_1)) &+ \omega \quad \mathcal{V}(M) \end{cases}$$

with ω = 0.499 (STs), ω = 0.048 (GTs), ω = 0.276 (IGs)

TDRL Model Dopamine D Applications

Meta-Learning slide # 111 / 141



with ω = 0.499 (STs), ω = 0.048 (GTs), ω = 0.276 (IGs)

slide # 112 / 141

Behavioral results



TDRL Model Dopamine D Applications

MMeta-Learning slide # 113 / 141

Physiological results





Physiological results

slide # 115 / 141

Pharmacological results



Dopamine TD Applications MMeta-Learning slide # 116 / 141

Summary of the simulation results



TDRL Model Dopamine TD Applications MMeta-Learning

<u>slide # 117 / 141</u>

Experimental predictions

- DA dip at each magazine visit during ITI.
- DA patterns in the intermediate group.
- Shortening the ITI should change DA pattern in GTs.
- Removing the magazine during ITI should abolish the difference in DA patterns between STs and GTs.
- Reducing the ITI duration should increase the tendency to goal-track in the overall population.

TDRL Model Dopamine D Applications Meta-Learning

MMeta-Learning slide # 118 / 141

META-LEARNING (regulation of decision-making) 1. Dual-system RL coordination

2. Online parameters tuning

META-LEARNING

TDRL Model Dopamine TD Applications MMeta-Learning slide # 119 / 141



META-LEARNING



TDRL Model Dopamine D Applications

Meta-Learning

META-LEARNING

• Meta-learning methods propose to tune RL parameters as a function of average reward and uncertainty (Schweighofer & Doya, 2003).

Meta-Learning

slide #



 \rightarrow Can we use such meta-learning principles to better understand neural mechanisms in the prefrontal cortex ?

Back to my post-doc work

Dopamine TD Applications MMeta-Learning

slide # 122 / 141



Question: How did the monkeys learn to re-explore after each presentation of the PCC signal?

Hypothesis: By trial-and-error during pretraining.

TDRL Model Dopamine D Applications

MMeta-Learning slide # 123 / 141



Khamassi et al. (2011) Frontiers in Neurorobotics

Robotic model of monkey behavior in this task

TDRL Model Dopamine D Applications

MMeta-Learning slide # 124 / 141

Dopamine TD Applications MMeta-Learning slide # 125 / 141

• Reproduction of the global properties of monkey performance in the PS task.



Khamassi et al. (2011) Frontiers in Neurorobotics



Model-based analysis My post-doc work

TDRL Model Dopamine D Applications

MMeta-Learning slide # 126 / 141



Multiple regression analysis with bootstrap

Khamassi et al. (2013) Prog Brain Res; Khamassi et al. (2014) Cerebral Cortex

Meta-learning applied to Human-Robot Interaction

TDRL Model Dopamine TD Applications

MMeta-Learning slide # 127 / 141



- In the previous task, monkeys and the model a priori 'know' that *PCC* means a reset of exploration rate and action values.
 - Here, we want the iCub robot to learn it by itself.

•

Application to simple learning in humanoid robot

TDRL Model Dopamine TD Applications

MMeta-Learning slide # 128 / 141



Khamassi et al. (2011) Frontiers in Neurorobotics

Meta-learning applied to Human-Robot Interaction

MMeta-Learning slide # 129 / 141



Error

Human's hands

Cheating

Cheating

Meta-learning applied to Human-Robot Interaction

TDRL Model Dopamine TD Applications

MMeta-Learning slide # 130 / 141

meta-value(i) \leftarrow meta-value(i) + α '. Δ [averageReward]



CONCLUSION OF THE ACC-LPFC META-LEARNING PART

TDRL Model Dopamine TD Applications MMeta-Learning slide # 131 / 141

- ACC is in an appropriate position to evaluate feedback
 history to modulate the exploration rate in LPFC.
- ACC-LPFC interactions could regulate exploration based on mechanisms capturable by the metalearning framework.
- Such modulation could be subserved via noradrenaline innervation in LPFC.
- Such a pluridisciplinary approach can contribute both to a better understanding of the brain and to the design of algorithms for autonomous decision-making.

Meta-learning and motor learning

Dopamine TD Applications MMeta-Learning

<u>slide # 132 / 141</u>

 Can meta-learning principles be useful for the integration of reinforcement learning and motor learning?

Structure learning (Braun Aertsen Wolpert Mehring 2009)



MMeta-Learning slide # 133 / 141

Structure learning (Braun Aertsen Wolpert Mehring 2009)

TDRL Model Dopamine TD Applications

MMeta-Learning slide # 134 / 141



Structure learning (Braun Aertsen Wolpert Mehring 2009)

MMeta-Learning slide # 135 / 141



Schmidhuber on meta-learning (1)

Recurrent neural-networks applied to Robotics



MMeta-Learning slide # 136 / 141

Mayer et al. (IROS 2006)

Schmidhuber on meta-learning (2)

Dopamine TD Applications MMeta-Learning slide # 137 / 141

- RL with self-modifying policies (actions that can edit the policy itself)
- Success-story criterion (time varying set V of past checkpoints that led to long-term reward accelerations)

Schmidhuber on motor learning

Dopamine TD Applications MMeta-Learning slide # 138 / 141

- Learning maps of task-relevant motor behaviors under specified constraints (e.g. maintain hands parallel; do not touch box nor table; ...)
- How can these primitive constrained motor behaviors be used by decision system and high-level goaldirected learning?



Stollenga et al. (IROS 2013)



- Dopamine neurons encode a reward prediction error.
- Model-based analysis in Neurosci of Decision-making
- Reinforcement Learning models need to be refined to explain behavior / neural activity:
 - multiple parallel decision systems.
 - off-line learning during sleep.
 - meta-learning (ACC-DLPFC interactions).
- These model improvements can produce testable experimental predictions (Pavlovian autoshaping ; Navigation ; L-DOPA in Parkinson disease ; …)

- The Reinforcement Learning framework provides algorithms for autonomous agents.
- It can also help explain neural activity in the brain.
- Such a pluridisciplinary approach can contribute both to a better understanding of the brain and to the design of algorithms for autonomous decision-making.

ACKNOWLEDGMENTS

TDRL Model Dopamine TD Applications

slide # 141 / 141

ISIR (CNRS – UPMC)

Nassim Aklil

Jean Bellot

Ken Caluwaerts

Raja Chatila

Laurent Dollé

Benoît Girard

Florian Lesaint

Olivier Sigaud

Guillaume Viejo

LAAS (CNRS Toulouse)

Rachid Alami Aurélie Clodic

Univ. Manchester

Mark D. Humphries

Financial support

FP6 IST 027189 European project



FINANCÉ PAR

Learning under Uncertainty Project ; ROBOERGOSUM Project



HABOT Project Emergence(s) Program