The Neuroscience of Reinforcement Learning

Yael Niv

Psychology Department & Neuroscience Institute Princeton University



ICML'09 Tutorial Montreal

 Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!
 - 1. Take pride

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!
 - 1. Take pride
 - 2. Ask: what can neuroscience do for me?

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!
 - 1. Take pride

2. Ask: what can neuroscience do for me?

• Why are you here?

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!
 - 1. Take pride

2. Ask: what can neuroscience do for me?

- Why are you here?
 - To learn about learning in animals and humans

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!
 - 1. Take pride

2. Ask: what can neuroscience do for me?

- Why are you here?
 - To learn about learning in animals and humans
 - To find out the latest about how the brain does RL

- Reinforcement learning has revolutionized our understanding of learning in the brain in the last 20 years
- Not many ML researchers know this!
 - 1. Take pride

2. Ask: what can neuroscience do for me?

- Why are you here?
 - To learn about learning in animals and humans
 - To find out the latest about how the brain does RL
 - To find out how understanding learning in the brain can help RL research

If you are here for other reasons...

learn what is RL and how to do it

learn about the brain in general

read email

take a wellneeded nap

smirk at the shoddy state of neuroscience

• The brain coarse-grain

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions



Credits: Daniel Wolpert 5



because computers were not yet invented

- because computers were not yet invented
- to behave

- because computers were not yet invented
- to behave

• example: sea squirt



- because computers were not yet invented
- to behave

• example: sea squirt



 larval stage: primitive brain & eye, swims around, attaches to a rock

- because computers were not yet invented
- to behave

• example: sea squirt



- larval stage: primitive brain & eye, swims around, attaches to a rock
- adult stage: sits. digests brain.







world









what do we know about the brain?

 Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)

what do we know about the brain?

 Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)



what do we know about the brain?

- Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)
- Single neurons: we know quite a bit about how they work (but still don't know much about how their 3D structure affects function)


- Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)
- Single neurons: we know quite a bit about how they work (but still don't know much about how their 3D structure affects function)





- Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)
- Single neurons: we know quite a bit about how they work (but still don't know much about how their 3D structure affects function)

- Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)
- Single neurons: we know quite a bit about how they work (but still don't know much about how their 3D structure affects function)
- Networks of neurons: we have some ideas but in general are still in the dark

- Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)
- Single neurons: we know quite a bit about how they work (but still don't know much about how their 3D structure affects function)
- Networks of neurons: we have some ideas but in general are still in the dark
- Learning: we know a lot of facts (LTP, LTD, STDP) (not clear which, if any are relevant; relationship between synaptic learning rules and computation essentially unknown)

- Anatomy: we know a lot about what is where and (more or less) which area is connected to which (But unfortunately names follow structure and not function; be careful of generalizations, e.g. neurons in motor cortex can respond to color)
- Single neurons: we know quite a bit about how they work (but still don't know much about how their 3D structure affects function)
- Networks of neurons: we have some ideas but in general are still in the dark
- Learning: we know a lot of facts (LTP, LTD, STDP) (not clear which, if any are relevant; relationship between synaptic learning rules and computation essentially unknown)
- Function: we have pretty coarse grain knowledge of what different brain areas do (mainly sensory and motor; unclear about higher cognitive areas; much emphasis on representation rather than computation)

Summary so far...

- We have a lot of facts about the brain
- But.. we still don't understand that much about how it works
- (can ML help??)

Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions

what do neuroscientists do all day?

what do neuroscientists do all day?

figure out how the brain generates behavior





Old idea:
 structure → function



Old idea:
 structure → function



- Old idea:
 structure → function
- The brain is an extremely complex (and messy) dynamic biological system



- Old idea:
 structure → function
- The brain is an extremely complex (and messy) dynamic biological system
- 10¹¹ neurons
 communicating through
 10¹⁴ synapses



- Old idea:
 structure → function
- The brain is an extremely complex (and messy) dynamic biological system
- 10¹¹ neurons
 communicating through
 10¹⁴ synapses
- we don't stand a chance...





• (relatively) New Idea:



- (relatively) New Idea:
- The brain is a computing device



- (relatively) New Idea:
- The brain is a computing device
- Computational models can help us talk about functions of the brain in a precise way



- (relatively) New Idea:
- The brain is a computing device
- Computational models can help us talk about functions of the brain in a precise way
- Abstract and formal theory can help us organize and interpret (concrete) data



David Marr (1945-1980) proposed three levels of analysis:

David Marr (1945-1980) proposed three levels of analysis:

1. the problem (Computational Level)

David Marr (1945-1980) proposed three levels of analysis:

- 1. the problem (Computational Level)
- 2. the strategy (Algorithmic Level)

David Marr (1945-1980) proposed three levels of analysis:

- 1. the problem (Computational Level)
- 2. the strategy (Algorithmic Level)
- 3. how its actually done by networks of neurons (Implementational Level)











SCHOOL IS H	ELL KING
SHOULD YOU	30
TO GRAD SCH	00L?
A WEE TEST	
I AM A COMPU	CSIVE
CRUSHED INTO D	UST.
PROFESSOR'S S	S A LAVE,
TIME IS USING AND CITING AUT	5000 JARGON IORITIES.
TO CONTINUE TO OF AVOIDING LIF	NEED te Process fe.

optimal decision making (maximize reward, minimize punishment)





SCHOOL IS HELL BUT IT BEATS WORKING
SHOULD YOU GO
TO GRAD SCHOOL?
A WEE TEST
I AM A COMPULSIVE
CRUSHED INTO DUST.
PROFESSOR'S SLAVE.
TIME IS USING JARGON AND CITING AUTHORITIES.
I FEEL A DEEP NEED TO CONTINUE THE PROCESS OF AVOIDING LIFE.

optimal decision making (maximize reward, minimize punishment)



Why is this hard?



SCHOOL IS HELL BUT IT BEATS WORKING
SHOULD YOU GO
TO GRAD SCHOOL?
A WEE TEST
I AM A COMPULSIVE
CRUSHED INTO DUST.
PROFESSOR'S SLAVE.
TIME IS USING JARGON AND CITING AUTHORITIES.
I FEEL A DEEP NEED TO CONTINUE THE PROCESS OF AVOIDING LIFE.

optimal decision making (maximize reward, minimize punishment)





		1. I. I. I.		
VVhv	IS	this	har	יםי
<u> </u>				

Reward/punishment may be delayed

SCHOOL IS HELL IT BEATS WORKING
SHOULD YOU GO
TO GRAD SCHOOL?
A WEE TEST
I AM A COMPULSIVE
CRUSHED INTO DUST.
PROFESSOR'S SLAVE.
MY IDEA OF A GOOD TIME IS USING JARGON
TO CONTINUE THE PROCESS OF AVOIDING LIFE.

optimal decision making (maximize reward, minimize punishment)





Why is this hard?

- Reward/punishment may be delayed
- Outcomes may depend on a series of actions

SCHOOL IS HELL
IT BEATS WORKING
SHOULD YOU GO
TO GRAD SCHOOL?
A WEE TEST
I AM A COMPULSIVE
CRUSHED INTO DUST.
PROFESSOR'S SLAVE.
TIME IS USING JARGON AND CITING AUTHORITIES.
I FEEL A DEEP NEED TO CONTINUE THE PROCESS OF AVOIDING LIFE.

optimal decision making (maximize reward, minimize punishment)





Why is this hard?

- Reward/punishment may be delayed
- Outcomes may depend on a series of actions
- ⇒ "credit assignment problem" (Sutton, 1978)

SCHOOL IS HELL BUT IT BEATS WORKING
SHOULD YOU GO
TO GRAD SCHOOL?
A WEE TEST
I AM A COMPULSIVE
CRUSHED INTO DUST.
PROFESSOR'S SLAVE.
TIME IS USING JARGON AND CITING AUTHORITIES.
TO CONTINUE THE PROCESS OF AVOIDING LIFE.

in comes reinforcement learning
in comes reinforcement learning

• The problem: optimal decision making (maximize reward, minimize punishment)

in comes reinforcement learning

- The problem: optimal decision making (maximize reward, minimize punishment)
- An algorithm: reinforcement learning

in comes reinforcement learning

- The problem: optimal decision making (maximize reward, minimize punishment)
- An algorithm: reinforcement learning
- Neural implementation: basal ganglia, dopamine

Summary so far...

- <u>Idea</u>: study the brain as a computing device
- Rather than look at what networks of neurons in the brain represent, look at what they compute
- What do animal's brains compute?

Animal Conditioning and RL

- two basic types of animal conditioning (animal learning)
- how do these relate to RL?









































Pavlovian conditioning examples (conditioned suppression, autoshaping)

Pavlovian conditioning examples (conditioned suppression, autoshaping)



how is this related to RL?



model-free learning of values of stimuli through experience; responding conditioned on (predictive) value of stimulus

Rescorla & Wagner (1972)

The idea: error-driven learning

Change in value is proportional to the difference between actual and predicted outcome

$$\Delta V(S_i) = \eta [R - \sum_{j \in \text{trial}} V(S_j)]$$

Two assumptions/hypotheses:

[1] learning is driven by error (formalize notion of surprise)[2] summations of predictors is linear

Phase I

Phase II

Phase I

Phase II



Phase I

Phase II





25

Phase I

Phase II





25





















• Background: Darwin, attempts to show that animals are intelligent

- Background: Darwin, attempts to show that animals are intelligent
- Thorndike (age 23): submitted
 PhD thesis on "Animal intelligence: an experimental study of the associative processes in animals"

- Background: Darwin, attempts to show that animals are intelligent
- Thorndike (age 23): submitted
 PhD thesis on "Animal intelligence: an experimental study of the associative processes in animals"
- Tested hungry cats in "puzzle boxes"





- Background: Darwin, attempts to show that animals are intelligent
- Thorndike (age 23): submitted
 PhD thesis on "Animal intelligence: an experimental study of the associative processes in animals"
- Tested hungry cats in "puzzle boxes"
- Definition for learning: time to escape





- Background: Darwin, attempts to show that animals are intelligent
- Thorndike (age 23): submitted
 PhD thesis on "Animal intelligence: an experimental study of the associative processes in animals"
- Tested hungry cats in "puzzle boxes"
- Definition for learning: time to escape
- Gradual learning curves, did not look like 'insight' but rather trial and error





Example: Free operant conditioning (Skinner)



Example: Free operant conditioning (Skinner)





how is this related to RL?



animals can learn an arbitrary policy to obtain rewards (and avoid punishments)

Summary so far...

- The world presents animals/humans with a huge reinforcement learning problem (or many such small problems)
- Optimal learning and behavior depend on prediction and control
- How can the brain realize these?
 Can RL help us understand the brain's computations?


Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions





Parkinson's Disease



→ Motor control / initiation?



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

- → Reward pathway?
- \rightarrow Learning?



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?

→ Learning?



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?

→ Learning?

Also involved in:

• Working memory



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?

→ Learning?

- Working memory
- Novel situations



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?

→ Learning?

- Working memory
- Novel situations
- ADHD



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?

→ Learning?

- Working memory
- Novel situations
- ADHD
- Schizophrenia



→ Motor control / initiation?

Drug addiction, gambling, Natural rewards

→ Reward pathway?

→ Learning?

Also involved in:

- Working memory
- Novel situations
- ADHD
- Schizophrenia

•••

 \bigcirc

role of dopamine: many hypotheses

- Anhedonia hypothesis
- Prediction error hypothesis
- Salience/attention
- (Uncertainty)
- Incentive salience
- Cost/benefit computation
- Energizing/motivating behavior

 Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus



- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia



- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning



- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning



- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning





- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning



- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning



- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning

- Anhedonia = inability to experience positive emotional states derived from obtaining a desired or biologically significant stimulus
- Neuroleptics (dopamine antagonists) cause anhedonia
- Dopamine is important for reward-mediated conditioning

but...



predictable reward





(Schultz et al. '90s) 35







Schultz, Dayan, Montague, 1997³⁶

looks familiar?



Schultz, Dayan, Montague, 1997³⁶



looks familiar?





prediction error hypothesis of dopamine

prediction error hypothesis of dopamine

The idea: Dopamine encodes a temporal difference reward prediction error (Montague, Dayan, Barto mid 90's)

prediction error hypothesis of dopamine



The idea: Dopamine encodes a temporal difference reward prediction error

(Montague, Dayan, Barto mid 90's)
prediction error hypothesis of dopamine



The idea: Dopamine encodes a temporal difference reward prediction error

(Montague, Dayan, Barto mid 90's)





37

prediction error hypothesis of dopamine



The idea: Dopamine p = 0.5encodes a temporal difference reward prediction error p = 0.75(Montague, Dayan, Barto mid 90's) -iorillo et al, 2003 p = 1.0stimulus on Tobler et al, 2005 0.0 ml 0.025 ml 0.075 ml 0.15 ml 0.25 ml ш 5 spikes/s

Onset of conditioned stimuli predicting expected reward value

37

reward

5 spikes

p = 0.0

p = 0.25

400





Bayer & Glimcher (2005) 38



$$V_{t} = \eta \sum_{i=1}^{t} (1 - \eta)^{t - i} r_{i}$$



Bayer & Glimcher (2005) 38





$$V_{t} = \eta \sum_{i=1}^{t} (1 - \eta)^{t-i} r_{i}$$



Bayer & Glimcher (2005) 38

where does dopamine project to?

main target: basal ganglia





• prediction errors are for learning...

- prediction errors are for learning...
- cortico-striatal synapses show dopamine-dependent plasticity

- prediction errors are for learning...
- cortico-striatal synapses show dopamine-dependent plasticity



- prediction errors are for learning...
- cortico-striatal synapses show dopamine-dependent plasticity





- prediction errors are for learning...
- cortico-striatal synapses show dopamine-dependent plasticity
- three-factor learning rule: need presynaptic+postsynaptic+dopamine





Conditioning can be viewed as prediction learning

• The problem: prediction of future reward

- The problem: prediction of future reward
- The algorithm: temporal difference learning

- The problem: prediction of future reward
- The algorithm: temporal difference learning
- Neural implementation: dopamine dependent learning in corticostriatal synapses in the basal ganglia

- The problem: prediction of future reward
- The algorithm: temporal difference learning
- Neural implementation: dopamine dependent learning in corticostriatal synapses in the basal ganglia
- ⇒ Precise (normative!) theory for generation of dopamine firing patterns

- The problem: prediction of future reward
- The algorithm: temporal difference learning
- Neural implementation: dopamine dependent learning in corticostriatal synapses in the basal ganglia
- ⇒ Precise (normative!) theory for generation of dopamine firing patterns
- ⇒ A computational model of learning allows us to look in the brain for "hidden variables" postulated by the model

Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions

3 model-free learning algorithms

Actor/Critic Q learning SARSA

Actor/Critic in the brain?



Actor/Critic in the brain?



Actor/Critic in the brain?



evidence for this?





 measure BOLD ("blood oxygenation level dependent") signal



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet

<u>ldea</u>:

• Brain is functionally modular



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet

- Brain is functionally modular
- Neural activity uses energy & oxygen



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet

- Brain is functionally modular
- Neural activity uses energy & oxygen
- Measure brain usage, not structure



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet

- Brain is functionally modular
- Neural activity uses energy & oxygen
- Measure brain usage, not structure
- Spatial resolution: ~3mm 3D "voxels"



- measure BOLD ("blood oxygenation level dependent") signal
- oxygenated vs de-oxygenated hemoglobin have different magnetic properties
- detected by big superconducting magnet

- Brain is functionally modular
- Neural activity uses energy & oxygen
- Measure brain usage, not structure
- Spatial resolution: ~3mm 3D "voxels"
- temporal resolution: 5-10 seconds


short aside: functional magnetic resonance imaging (fMRI)



short aside: functional magnetic resonance imaging (fMRI)



















• cond 1: instrumental (choose stimuli) - show preference for high probability stimulus in rewarding but not neutral trials



- cond 1: instrumental (choose stimuli) show preference for high probability stimulus in rewarding but not neutral trials
- cond 2: Pavlovian only indicate the side the 'computer' has selected (RTs as measure of learning)



- cond 1: instrumental (choose stimuli) show preference for high probability stimulus in rewarding but not neutral trials
- cond 2: Pavlovian only indicate the side the 'computer' has selected (RTs as measure of learning)
- why was the experiment designed this way (hint: think of prediction errors)

ventral striatum: correlated with prediction error in both conditions





B Instrumental







ventral striatum: correlated with prediction error in both conditions







Dorsal striatum: prediction error only in instrumental task



do prediction errors really influence learning?



Schonberg et al. 2007 49

do prediction errors really influence learning?





Schonberg et al. 2007 ⁵⁰

Summary so far...

- Some evidence for an Actor/Critic architecture in the brain
- Links predictions (Critic) to control (Actor) in very specific way; assumes no Q values
- (Not at all conclusive evidence)



Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions



Morris et al. 2005 53



Morris et al. 2005 ⁵⁴

stimulus on



stimulus on





















Differences from Morris et al. (2005):



Differences from Morris et al. (2005): • rats not monkeys



Differences from Morris et al. (2005): • rats not monkeys • VTA not SNc



Differences fromMorris et al. (2005):rats not monkeysVTA not SNc

• amount of training



Differences from Morris et al. (2005):

- rats not monkeys
- VTA not SNc
- amount of training
- task (representation of stimuli?)



Differences from Morris et al. (2005):

- rats not monkeys
- VTA not SNc
- amount of training
- task (representation of stimuli?)

(notice the messy signal... due to measurement or is it that way in the brain?)

Summary so far...

- SARSA or Q-learning? The jury is still out
- What needs to be done: more experiments recording from dopamine in telltale tasks
- The brain (dopamine) can inform RL: how does it learn in real time, with real noise, in real problems?
Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions

training:



training:





training:





training:





Even the humble rat can can learn spatial structure, and use it to plan flexibly

1 - Training:



1 - Training:



2 – Pairing with illness:





1070°¢



2 – Pairing with illness:

2 - Motivational shift:







1 - Training:

2 – Pairing with illness:

2 – Motivational shift:





Hungry

0 0





3 - Test:
(no rewards)



61











Animals will sometimes work for food they don't want!



Animals will sometimes work for food they don't want!



Animals will sometimes work for food they don't want!
 → in daily life: actions become automatic with repetition



Animals will sometimes work for food they don't want!
 → in daily life: actions become automatic with repetition



Animals will sometimes work for food they don't want!
 → in daily life: actions become automatic with repetition

overtrained rats



overtrained rats



overtrained rats



Animals with lesions to DLS never develop habits despite extensive training

Yin et al (2004), Coutureau & Killcross (2003) 63

overtrained rats



- → animals with lesions to DLS never develop habits despite extensive training
- → also treatments depleting dopamine in DLS

overtrained rats



- → animals with lesions to DLS never develop habits despite extensive training
- → also treatments depleting dopamine in DLS
- → also inactivations of infralimbic PFC after training



Yin, Ostlund, Knowlton & Balleine (2005) 64



lesions of the pDMS cause animals to leverpress habitually even with only moderate training



lesions of the pDMS cause animals to leverpress habitually even with only moderate training (also.. pre-limbic PFC, dorsomedial thalamus)

• The same action (leverpressing) can arise from two psychologically dissociable pathways

- The same action (leverpressing) can arise from two psychologically dissociable pathways
 - 1. moderately trained behavior is "goal-directed": dependent on outcome representation

- The same action (leverpressing) can arise from two psychologically dissociable pathways
 - 1. moderately trained behavior is "goal-directed": dependent on outcome representation
 - 2. overtrained behavior is "habitual": apparently not dependent on outcome representation

- The same action (leverpressing) can arise from two psychologically dissociable pathways
 - 1. moderately trained behavior is "goal-directed": dependent on outcome representation
 - 2. overtrained behavior is "habitual": apparently not dependent on outcome representation
- Lesions suggest two parallel systems; the intact one can apparently support behavior at any stage

- The same action (leverpressing) can arise from two psychologically dissociable pathways
 - 1. moderately trained behavior is "goal-directed": dependent on outcome representation
 - 2. overtrained behavior is "habitual": apparently not dependent on outcome representation
- Lesions suggest two parallel systems; the intact one can apparently support behavior at any stage
- Can RL help us make sense of this mess?

strategy 1: model based RL








learn model of task through experience





learn model of task through experience compute Q values by dynamic programming (or other method of lookahead/planning)





learn model of task through experience compute Q values by dynamic programming (or other method of lookahead/planning)

computationally costly, but also flexible (immediately sensitive to change)





learn model of task through experience compute Q values by dynamic programming (or other method of lookahead/planning)

computationally costly, but also flexible (immediately sensitive to change)







Stored:

$$Q(S_0,L) = 4$$

 $Q(S_0,R) = 2$
 $Q(S_1,L) = 4$
 $Q(S_1,R) = 0$
 $Q(S_1,R) = 0$
 $Q(S_1,R) = 0$
 $Q(S_1,R) = 0$



• learn values through prediction errors

Stored:

$$Q(S_0,L) = 4$$

 $Q(S_0,R) = 2$
 $Q(S_1,L) = 4$
 $Q(S_1,R) = 0$
 $Q(S_1,R) = 0$
 $Q(S_1,R) = 0$



- learn values through prediction errors
- choosing actions is easy so behavior is quick, reflexive

Stored:	$Q(S_1,L) =$
$Q(S_0,L) = 4$	$Q(S_1,R) =$
$Q(S_0, R) = 2$	$Q(S_2,L) =$
	$Q(S_2, R) =$



- learn values through prediction errors
- choosing actions is easy so behavior is quick, reflexive
- but needs a lot of experience to learn

Stored: $Q(S_0,L) = 4$ $Q(S_0,R) = 2$ $Q(S_1,L) = 4$ $Q(S_1,R) = 0$ $Q(S_1,R) = 0$ $Q(S_1,R) = 0$ $Q(S_1,R) = 0$



- learn values through prediction errors
- choosing actions is easy so behavior is quick, reflexive
- but needs a lot of experience to learn
- and inflexible, need relearning to adapt to any change (habitual)

Stored: $Q(S_0,L) = 4$ $Q(S_0,R) = 2$ $Q(S_1,L) = 4$ $Q(S_1,R) = 0$ $Q(S_1,R) = 0$ $Q(S_1,R) = 0$ $Q(S_1,R) = 0$

this answer raises two questions:



this answer raises two questions:

• Why should the brain use two different strategies/ controllers in parallel?



this answer raises two questions:

- Why should the brain use two different strategies/ controllers in parallel?
- If it uses two: how can it arbitrate between the two when they disagree (new decision making problem...)



1. each system is best in different situations (use each one when it is most suitable/most accurate)

- 1. each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) good with limited training, close to the reward (don't have to search ahead too far)

- 1. each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) good with limited training, close to the reward (don't have to search ahead too far)
 - habitual (cache) good after much experience, distance from reward not so important

- 1. each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) good with limited training, close to the reward (don't have to search ahead too far)
 - habitual (cache) good after much experience, distance from reward not so important
- 2. arbitration: trust the system that is more confident in its recommendation

- 1. each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) good with limited training, close to the reward (don't have to search ahead too far)
 - habitual (cache) good after much experience, distance from reward not so important
- 2. arbitration: trust the system that is more confident in its recommendation



- 1. each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) good with limited training, close to the reward (don't have to search ahead too far)
 - habitual (cache) good after much experience, distance from reward not so important
- 2. arbitration: trust the system that is more confident in its recommendation
 - use Bayesian RL (explore/exploit in unknown MDP; POMDP)



- 1. each system is best in different situations (use each one when it is most suitable/most accurate)
 - goal-directed (forward search) good with limited training, close to the reward (don't have to search ahead too far)
 - habitual (cache) good after much experience, distance from reward not so important
- 2. arbitration: trust the system that is more confident in its recommendation
 - use Bayesian RL (explore/exploit in unknown MDP; POMDP)
 - different sources of uncertainty in the two systems



Summary so far...

- animal conditioned behavior is not a simple unitary phenomenon: the same response can result from different neural and computational origins
- different neural mechanisms work in parallel to support behavior: cooperation and competition
- RL provides clues as to why this should be so, and what each system does

Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine (skipped for lack of time)
- Risk sensitivity and RL in the brain *NEW*
- Open challenges and future directions





sure \$20





or

sure \$20



risky \$40/\$0







1. Decision making is sensitive to risk





- 1. Decision making is sensitive to risk
- 2. RL (expected) values ignore risk





- 1. Decision making is sensitive to risk
- 2. RL (expected) values ignore risk
- 3. BOLD signals in nucleus accumbens correlate with prediction errors (can infer value of options from these)





- 1. Decision making is sensitive to risk
- 2. RL (expected) values ignore risk
- 3. BOLD signals in nucleus accumbens correlate with prediction errors (can infer value of options from these)

Will the neural value of a <u>sure 20¢</u> be the same as that of a 50% chance <u>risky 40¢</u> ?

Niv, Edlund, Dayan, O'Doherty (not yet published) 72



<u>same</u> neural values <u>different</u> neural values



<u>same</u> neural values <u>different</u> neural values

risk-sensitive decisions due to other mechanism that tracks risk

> "decision value" = E(r) + αV(r)

Kuhnen & Knutson(2005) Preuschoff, Bossaerts & Quartz (2006)



<u>same</u> neural values <u>different</u> neural values

risk-sensitive decisions due to other mechanism that tracks risk

73





<u>same</u> neural values

<u>different</u> neural values

risk-sensitive decisions due to other mechanism that tracks risk

risk-sensitive decisions due to sampling biases












































sure 20¢ vs. risky 40¢/0¢

<u>same</u> neural values <u>different</u> neural values
















































Sorry...

- This is of yet unpublished material so I hesitate to put the rest of this study online
- If you would like a copy of these slides or of the paper, feel free to email me at <u>yael@princeton.edu</u>

Although we are used to thinking about expected rewards in RL...

- Although we are used to thinking about expected rewards in RL...
- The brain (and human behavior) seems to fold risk (variance) into predictive values as well

- Although we are used to thinking about expected rewards in RL...
- The brain (and human behavior) seems to fold risk (variance) into predictive values as well
- Why is this a good thing to do?

- Although we are used to thinking about expected rewards in RL...
- The brain (and human behavior) seems to fold risk (variance) into predictive values as well
- Why is this a good thing to do?
- Can this help RL applications?

Outline

- The brain coarse-grain
- Learning and decision making in animals and humans: what does RL have to do with it?
- A success story: Dopamine and prediction errors
- Actor/Critic architecture in basal ganglia
- SARSA vs Q-learning: can the brain teach us about ML?
- Model free and model based RL in the brain
- Average reward RL & tonic dopamine
- Risk sensitivity and RL in the brain
- Open challenges and future directions

• How can RL deal with noisy inputs?

 How can RL deal with noisy inputs? 	 How does the brain deal with noisy inputs? (temporal noise!)

- How can RL deal with noisy inputs?
- How can RL deal with an unspecified state space?
- How does the brain deal with noisy inputs?
 (temporal noise!)

- How can RL deal with noisy inputs?
- How can RL deal with an unspecified state space?
- How does the brain deal with noisy inputs?
 (temporal noise!)
- How does the brain deal with an unspecified state space?

- How can RL deal with noisy inputs?
- How can RL deal with an unspecified state space?
- How can RL deal with multiple goals? Transfer between tasks?

- How does the brain deal with noisy inputs?
 (temporal noise!)
- How does the brain deal with an unspecified state space?

- How can RL deal with noisy inputs?
- How can RL deal with an unspecified state space?
- How can RL deal with multiple goals? Transfer between tasks?

- How does the brain deal with noisy inputs?
 (temporal noise!)
- How does the brain deal with an unspecified state space?
- How does the brain deal with multiple goals?
 Transfer between tasks?

- How can RL deal with noisy inputs?
- How can RL deal with an unspecified state space?
- How can RL deal with multiple goals? Transfer between tasks?

- How does the brain deal with noisy inputs?
 (temporal noise!)
- How does the brain deal with an unspecified state space?
- How does the brain deal with multiple goals?
 Transfer between tasks?

Summary: What have we learned here?

- RL has revolutionized how we think about learning in the brain
- Theoretical, but also practical (even clinical?) implications for neuroscience
- Neuroscience continues to be a "consumer" of ML theory/algorithms
- This does not have to be a one-way street: humans solve some problems so well that it is silly not to use human learning as an inspiration for new RL methods

THANK YOU!



interested in reading more? some recent reviews of neural RL

- Y Niv (2009) Reinforcement learning in the brain The Journal of Mathematical Psychology
- P Dayan & Y Niv (2008) Reinforcement learning and the brain: The Good, The Bad and The Ugly - Current Opinion in Neurobiology, 18(2), 185-196
- MM Botvinick, Y Niv & A Barto (2008) Hierarchically organized behavior and its neural foundations: A reinforcement learning perspective Cognition (online prepublication)
- K Doya (2008) Modulators of decision making Nature Neuroscience 11,410-416
- MFS Rushworth & TEJ Behrens (2008) Choice, uncertainty and value in prefrontal and cingulate cortex - Nature Neuroscience 11, 389-397
- A Johnson, MA van der Meer & AD Redish (2007) Integrating hippocampus and striatum in decision-making - Current Opinion in Neurobiology, 17, 692-697
- JP O'Doherty, A Hampton & H Kim (2007) Model-based fMRI and its application to reward learning and decision making Annals of the New York Academy of Science, 1104, 35-53
- ND Daw & K Doya (2006) The computational neurobiology of learning and reward -Current Opinion in Neurobiology, 6, 199-204