

# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

### Objectives for today:

- basic idea of policy gradient: learn actions, not Q-values
- log-likelihood trick: getting the correct statistical weight
- policy gradient algorithms
- why subtract the mean reward?
- reinforce with baseline (see actor critic)

## **Reading for this week:**

**Sutton and Barto, Reinforcement Learning  
(MIT Press, 2<sup>nd</sup> edition 2018, also online)**

Chapter: 13.1-13.5

**Background reading: none**

# **Announcements:**

Exam: Monday July ?, 12h15-15h15 (180 min)

## **Exam (counts 70%)**

- paper and pencil
- bring A5 sheet, handwritten notes, double-sided
- no textbook, no slides, no calculator.
- similar to exercises and quizzes
- sample exams on Moodle (from two last years)
  - smaller part: quiz-like questions
  - bigger part: exercise-like calculations
    - (typically 4 exercises with a, b, c, d ...)
- some points are 'easy', some medium, some 'difficult'

## **Miniproject (counts 30%).**

- miniproject validated after 'fraud detection interview'
- sign up for May 28<sup>th</sup> or June 4<sup>th</sup>

## **Recommended exam preparation**

- (1) do (or redo) **exercises** yourself
- (2) if stuck, read the relevant chapter of the **textbook**  
(see page 2 of slides of each week)
- (3) check the solution of exercise
- (4) look at the **quiz question** (always orange slides)
- (5) if stuck, read the relevant chapter of the **textbook**  
(see page 2 of slides of each week)
- (6) Look at **past exams** (solutions: see analog exercises)

**NOTE:** the slides are most useful if you have followed and annotated them yourself during the lecture; they are not designed as a stand-alone tool.

This is what successful students said about exam preparation:

Student A:

“For me, going through the exercises was very helpful, along with the slide quizzes. We also discussed theoretical questions from the lectures with my teammate and friends”

Student B:

“During the semester I have read the commented version of the slides in order to carry out the 2 miniprojects. I took care to understand each remark and I did the exercises when I had trouble in learning a topic. Before the exam, I felt that I was remembering well so I could focus only on Reinforcement Learning. In this case I found more useful solving the exercises to understand some key differences between the different algorithms e.g. off-policy versus on-policy.”

# This is what successful students said about exam preparation:

## Student C:

« I first went through all the lecture slides which I had taken notes on during lectures to reinforce my memory of various notions introduced in this course, and I want to stress that the comment pages were truly helpful. Afterwards, I went over all the exercises and collected a few questions to pose in the revision session held by TAs and got satisfactory clarification for most of them. »

## Student D:

« I prepared for the exam by reading slides over and over again. I think the comments slides helped me a lot in understanding and reading them over again helped me to build the structure of the overall course.

Exercises helped as well since it turns out that the exam is quite similar to exercises. »

# This is what successful students said about exam preparation:

## Student E:

“I attended nearly every class and made sure I understood the blackboard proofs properly because these were usually very useful for understanding the main concepts. During the exam preparation, I mostly just went through the class slides again and solved all of the exercises.”

## Student F:

“I never came to class but I did all the exercises and studied the books on Reinforcement Learning and Deep Learning.”

# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

### Objectives for today:

- basic idea of policy gradient: learn actions, not Q-values
- log-likelihood trick: getting the correct statistical weight
- policy gradient algorithms
- why subtract the mean reward?
- reinforce with baseline (see actor critic)



# 1. Review: Artificial Neural Networks for action learning



Where is the supervisor?  
Where is the labeled data?

Replaced by:

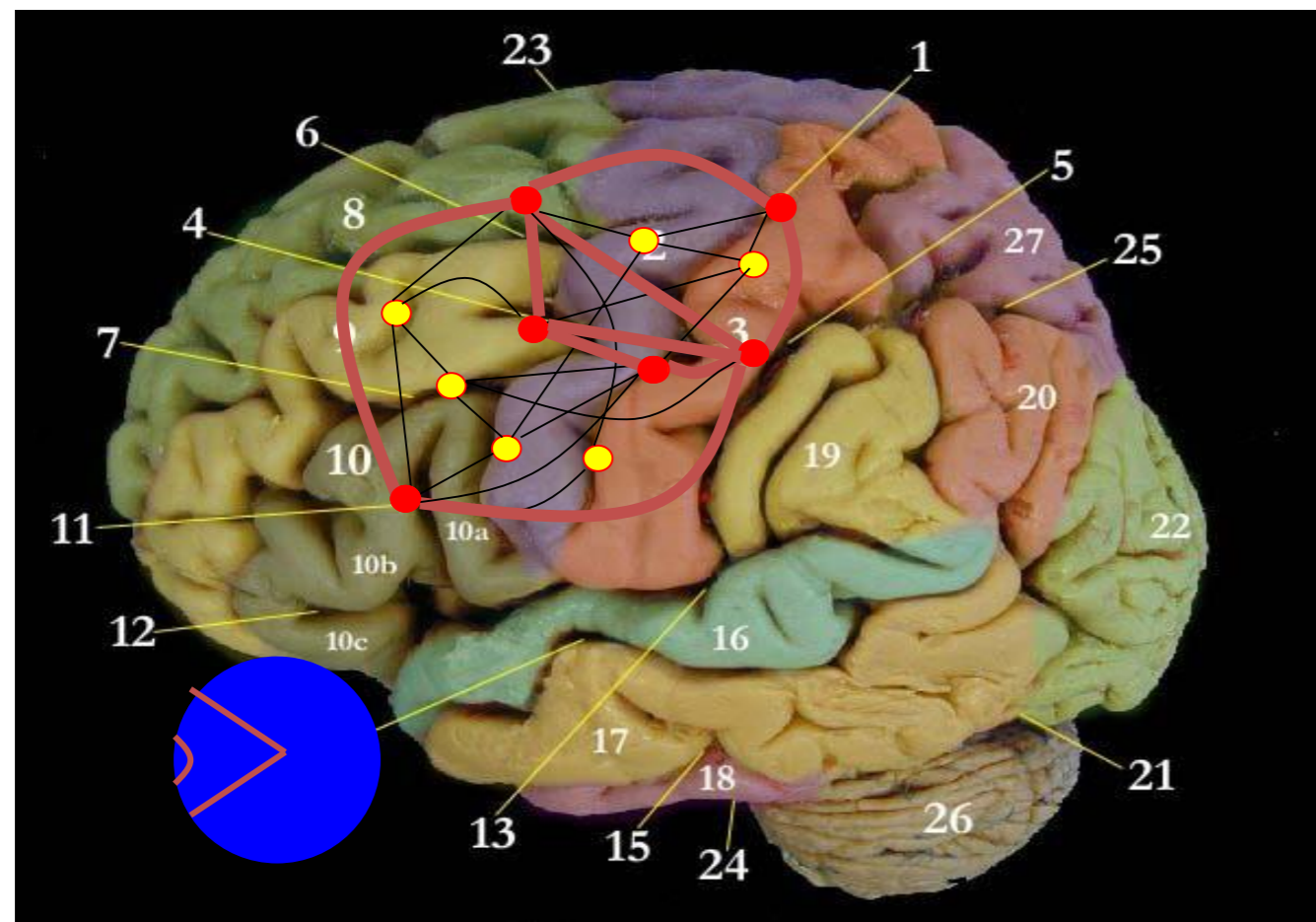
‘Value of action’

- ‘goodie’ for dog
- ‘success’
- ‘compliment’

BUT:

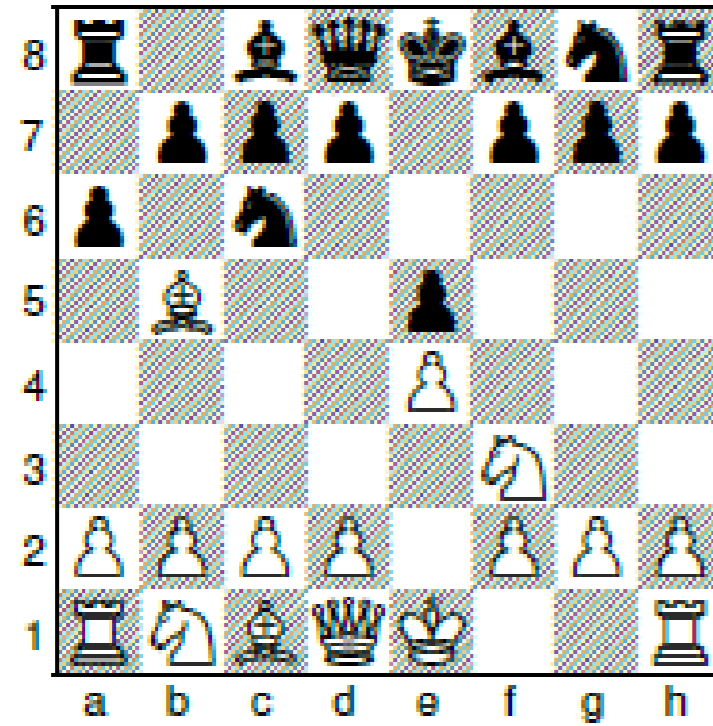
Reward is rare:

‘sparse feedback’ after  
a long action sequence



# 1. First steps toward Deep reinforcement learning

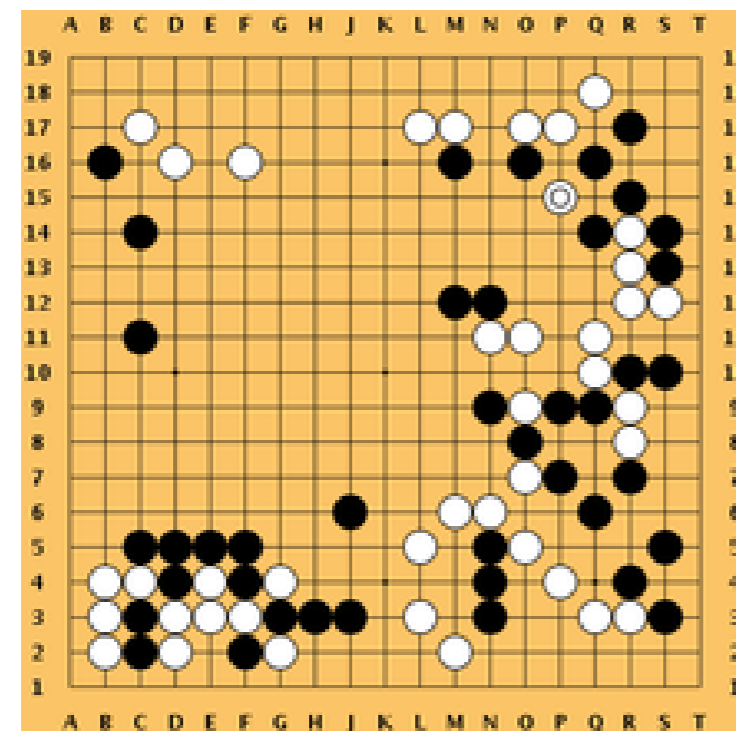
## Chess



Artificial neural network  
(*AlphaZero*) discovers different  
strategies by playing against itself.

In Go, it beats Lee Sedol

## Go

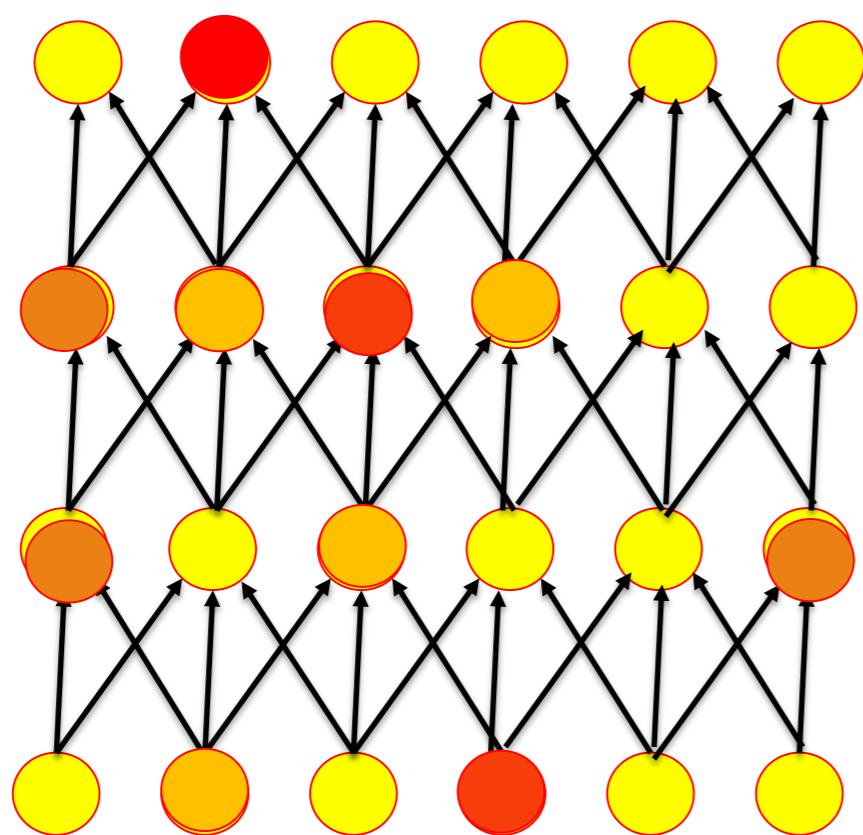


# 1. Backprop for deep Q-learning

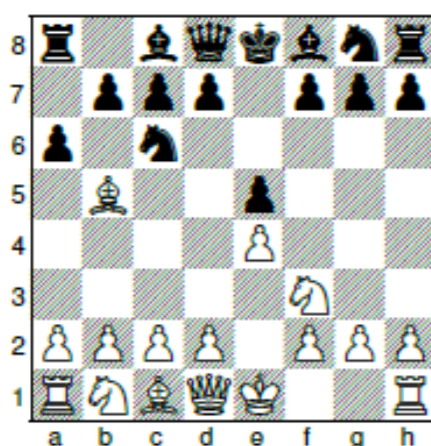
action and Q-values:

*Advance king*

output ↑ ↑ ↑ ↑ ↑



input ↑ ↑ ↑ ↑ ↑



Outputs are Q-values  
→ actions are easy to choose

For example:

Softmax strategy: take **action a'**  
with prob

$$P(a') = \frac{\exp[\beta Q(a')]}{\sum_a \exp[\beta Q(a)]}$$

(previous slide)

Last week we have seen that we can model Q-values in continuous state space as a function of the state  $s$ , and parameterized with weights  $w$ .

But in fact, a model of Q-values also works when the input space is discrete, such as it is in chess. Suppose that each output corresponds to one action (e.g. one type of move in chess).

We can use a neural network where the output are the Q-values of the different actions while the input represents the current state  $s$ .

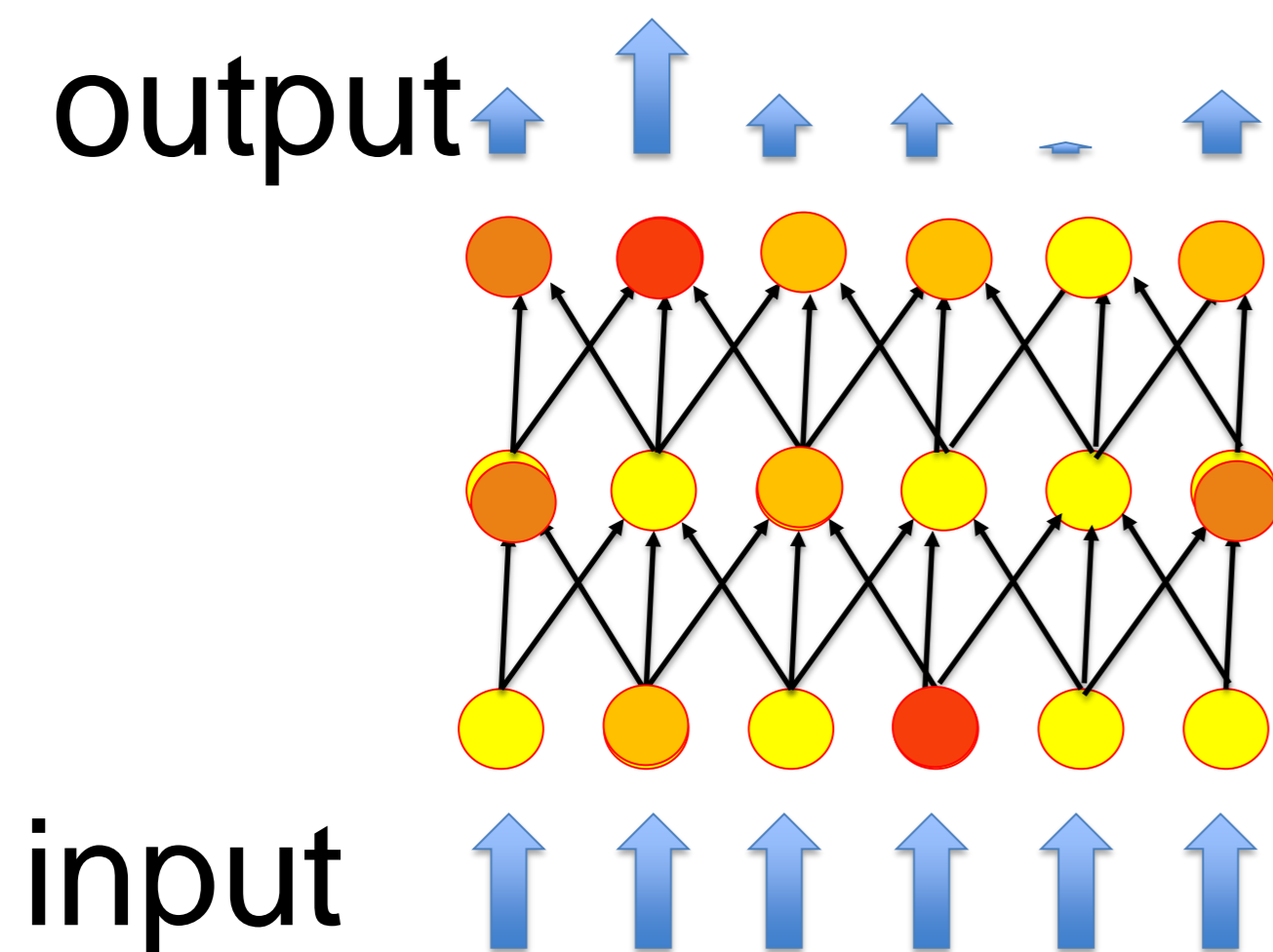
Thus, an output unit  $n$  represents  $Q(a_n, s)$ .

# 1. Backprop for deep Q-learning

(Backprop = gradient descent rule in multilayer networks)

action and Q-values:

*Move piece*



**Neural network parameterizes Q-values as a function of continuous state  $s$ .  
One output for one action  $a$ .  
Learn weights by playing against itself.**

Error function for SARSA

$$E = 0.5 [ r + \gamma Q(s',a') - Q(s,a) ]^2$$

(previous slide)

Suppose that each output corresponds to one action (e.g. one type of move in chess). Parameters are now the weights of the artificial neural network.

Actions are chosen, for example, by softmax on the Q-values in the output.

Weights are learned by playing against itself – doing gradient descent on an error function  $E$ .

Last week we finished by stating the error function:

$$E = 0.5 [ r + \gamma Q(s',a') - Q(s,a) ]^2$$

This error function will depend on the weights  $w$  (since  $Q(s,a)$  depends on  $w$ ). We can change the weights by gradient descent on the error function. This leads to the Backpropagation algorithm of 'Deep learning' (will be discussed next week).

# 1. Error function for continuous input representation

Consistency condition of Bellman Eq.

$$Q(s, a) = \sum_{s'} P_{s \rightarrow s'}^a \left[ R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi(s', a') Q(s', a') \right]$$

On-line consistency condition  
(should hold on average)

$$Q(s, a) = r + \gamma Q(s', a')$$

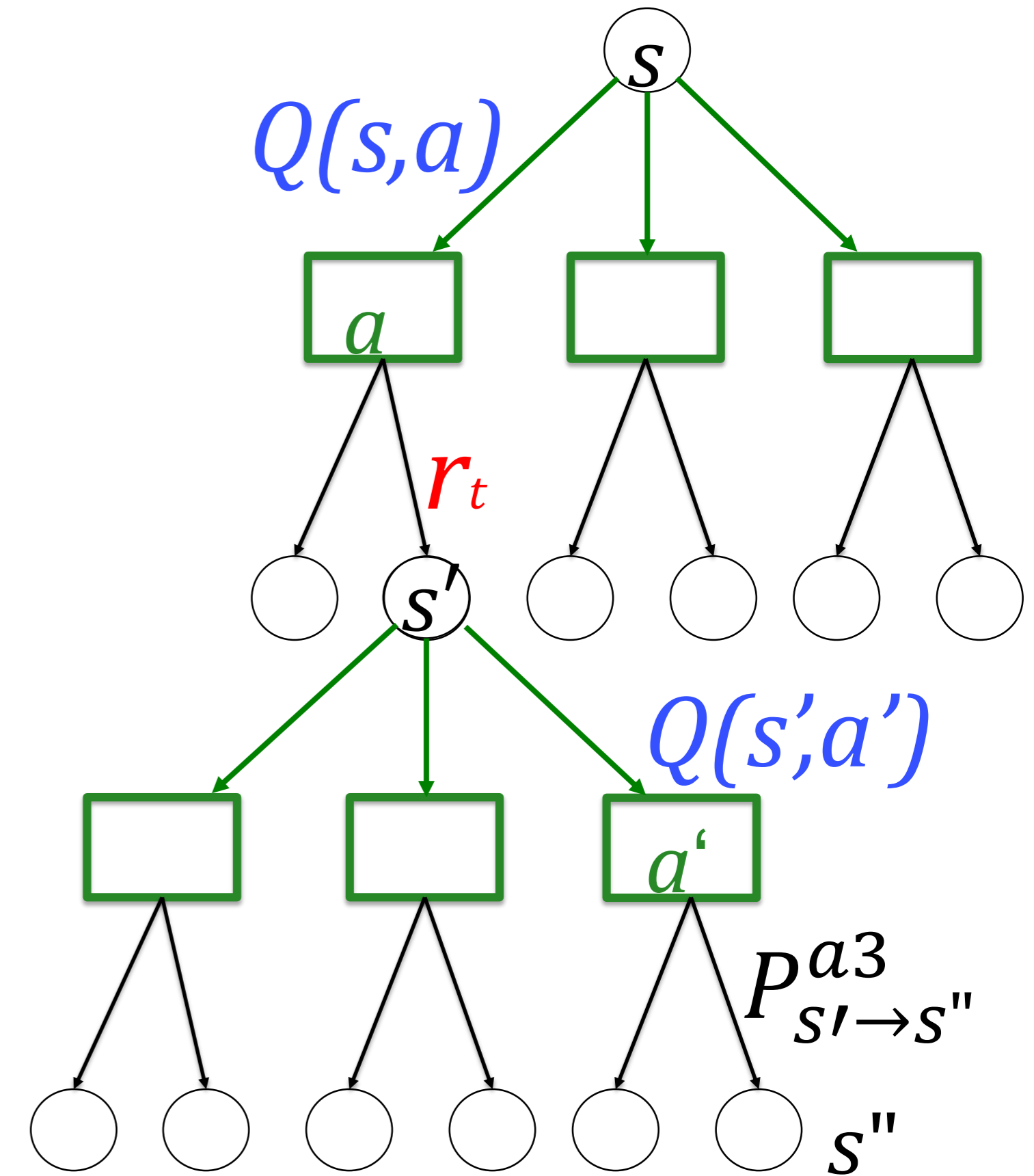
yields (online) Error function

target

$$E(\mathbf{w}) = \frac{1}{2} [r_t + \gamma Q(s', a' | \mathbf{w}) - Q(s, a | \mathbf{w})]^2$$

ignore

take gradient w.r.t  $\mathbf{w}$



(previous slide)

During the discussion of the Bellman equation and SARSA, we stated repeatedly that essentially we formulate a consistency condition.

$$Q(s,a) = r + \gamma Q(s',a')$$

Where the equality sign has to be interpreted as ‘should ideally on average be close to’ and the right hand side is the ‘target of learning’

The quadratic error function  $E$  measures how close we are to such an ideal case. This error function works not just for continuous state space, but also for a discrete state space such as in chess.

**IMPORTANT NOTE:** Since the ‘target of learning’ should be considered as momentarily fixed, we optimize the error function by taking the derivative of  $E$  with respect to the  $w$  in  $Q(s,a)$  but ignore that the target  $Q(s',a')$  also depends on  $w$ . In other words, during the optimization step we consider  $Q(s',a')$  as fixed.

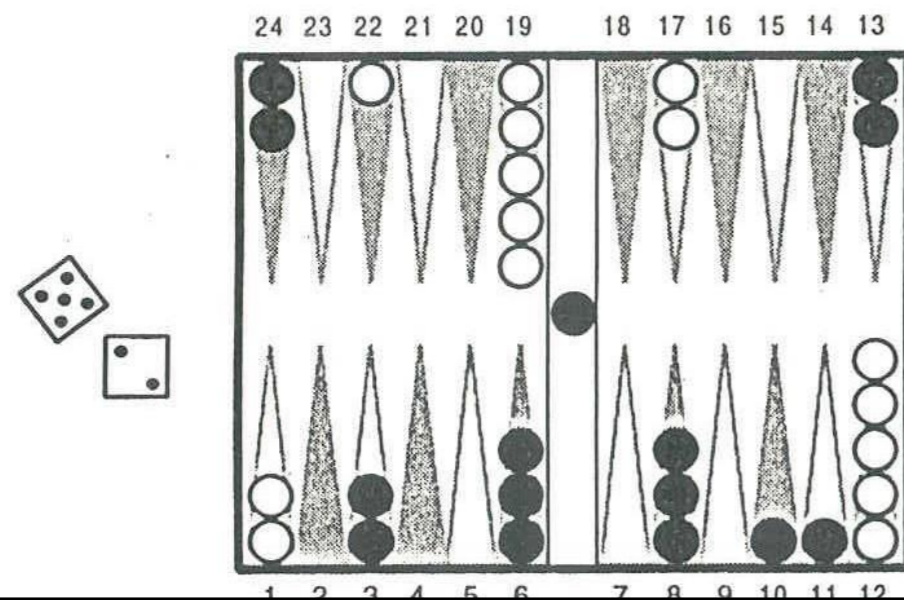
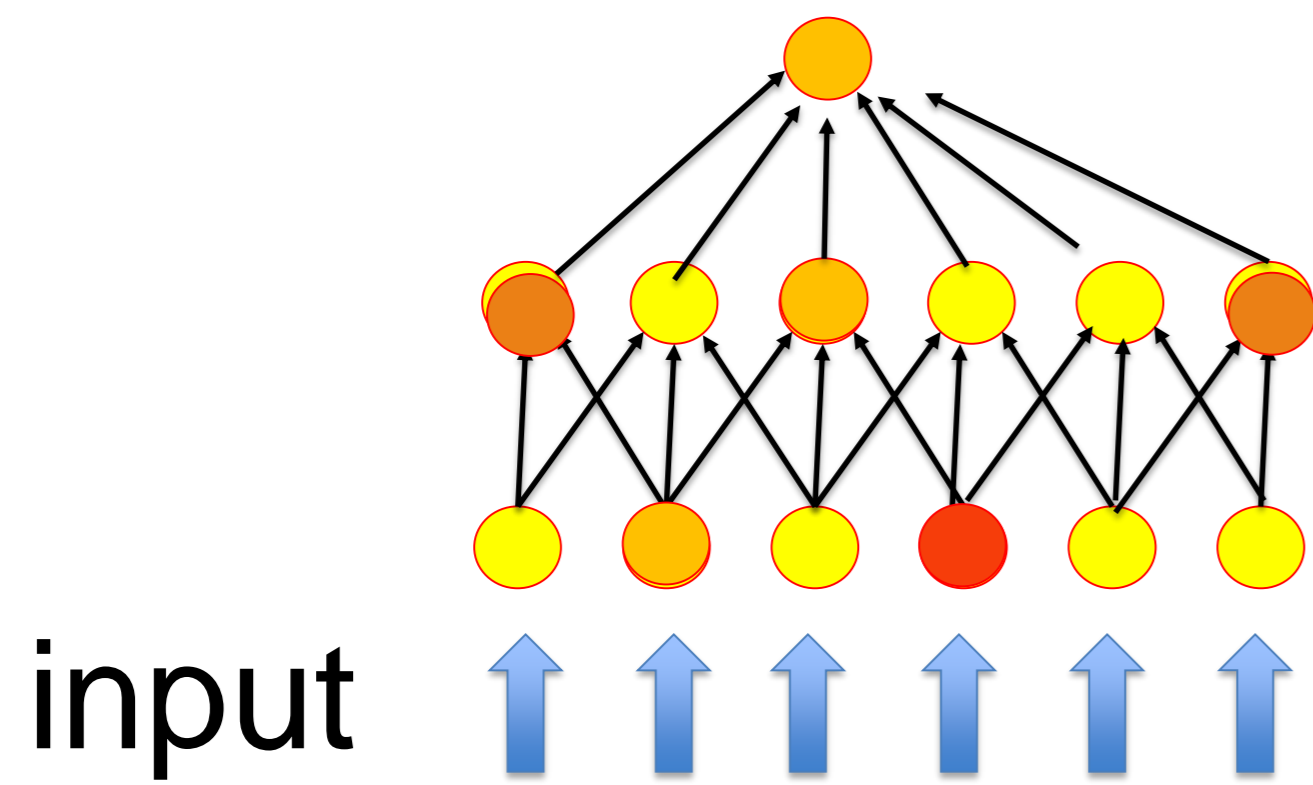
Taking the derivative with fixed target yields SARSA; see Exercise 5 of RL2.



# 1. Deep Neural Network for Value function

**Action:** move piece by epsilon greedy so as to increase V-value in each step

output: V-values:



- **Neural network parameterizes V-values as a function of state s.**
- **One single output.**
- **Learn weights by playing against itself.**
- **Minimize TD-error of V-function**
- **use eligibility traces**

TD-Gammon

Tesauro, 1992, 1994, **1995**, 2002

(previous slide)

The very same ideas can also be applied to learning the V-values, instead of Q-values. The advantage is that we have one single output. The disadvantage is that we need to look ahead (next possible states) to choose the action. But for games with a small number of 'possible next states' this is not a problem.

The analogous Bellman equation for the V-values leads to a consistency condition characterized by an error function

$$E(\mathbf{w}) = \frac{1}{2} [r_t + \gamma V(s'|\mathbf{w}) - V(s|\mathbf{w})]^2$$

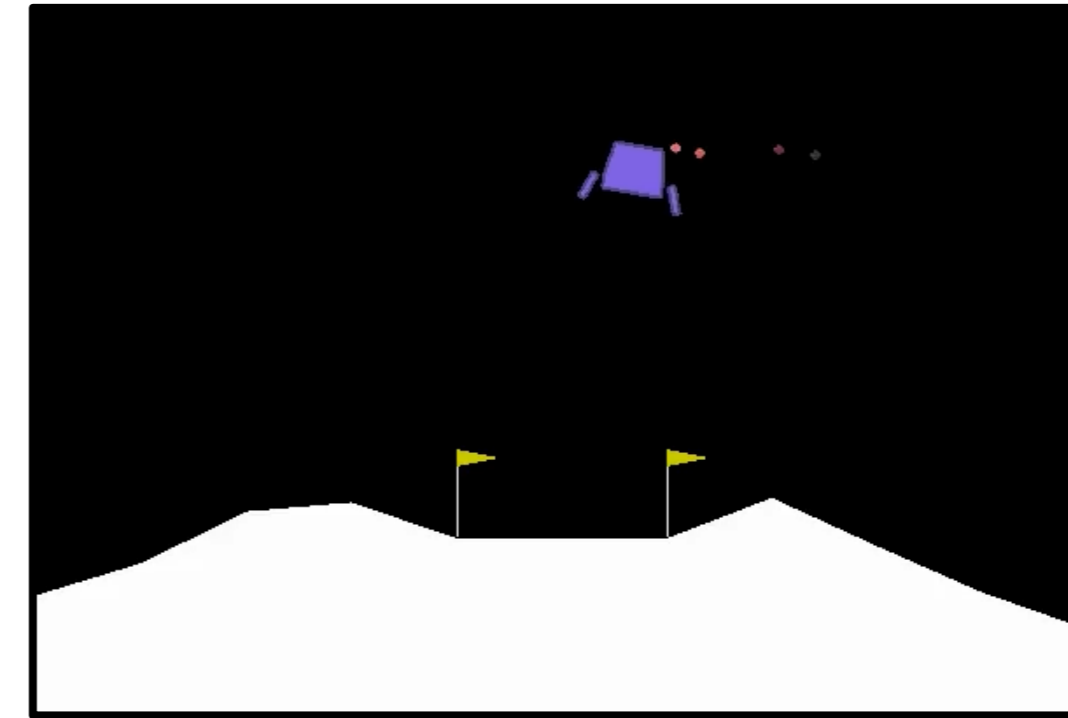
Eligibility traces enable to connect the reward at the end to states several steps before.

# 1. Neural networks to model input space

- for control problems, input space is naturally continuous

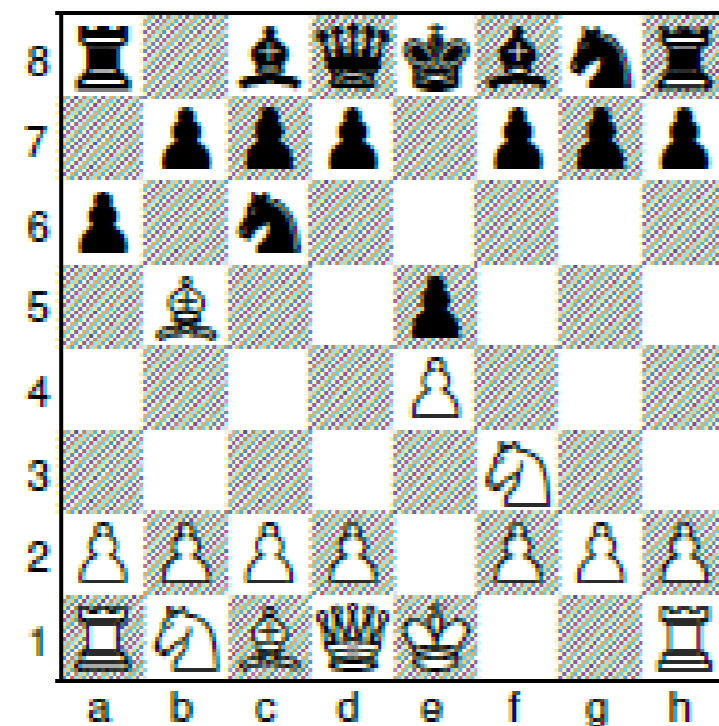
Moon lander

Aim: land between poles

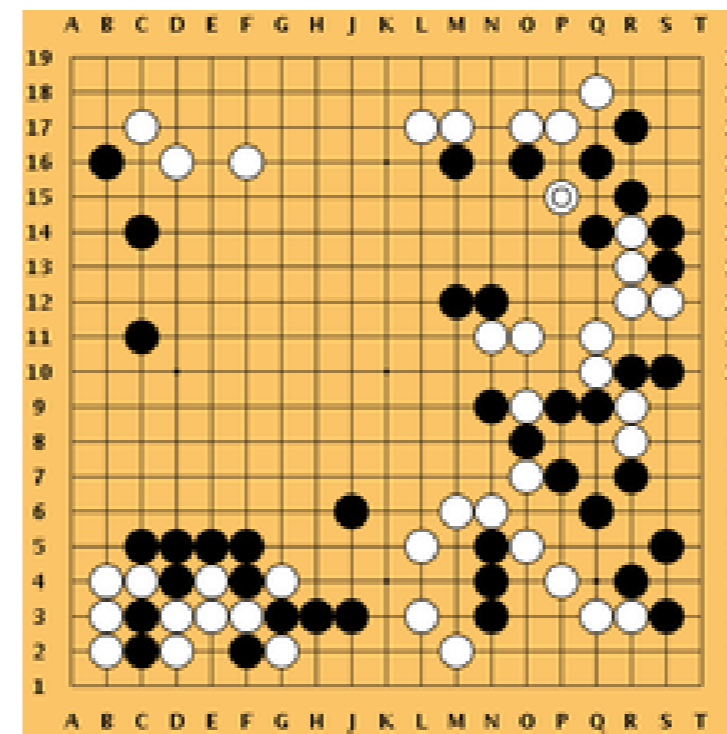


- for discrete games, the input space often too big  
→ generalize via hidden states in neural networks

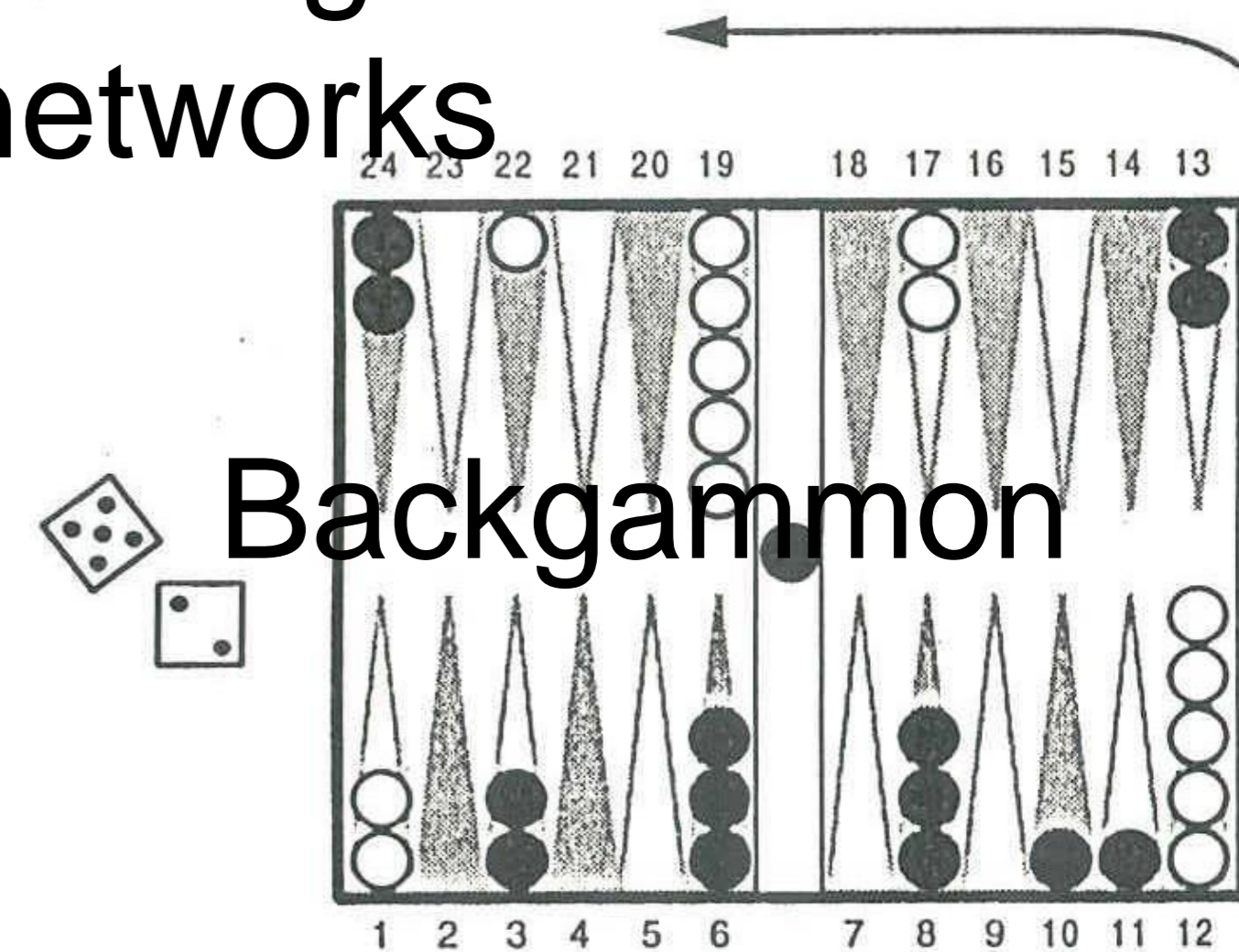
Chess



Go



Backgammon



(previous slide)

Why is it useful to use a continuous (as apposed to tabular) description of input space even in cases where the input is naturally discrete such as in games?

The reason is that describing Q-values as a SMOOTH function of the input enables generalization. Hidden layers of neural networks are able to extract compressed representations of the input space that introduce heuristic but useful notion of what it means that two states are 'similar' or 'neighbors'.

Related ideas have been used in many other applications, beyond chess backgammon or Go. We will study some of these later in this class.

TD learning where Q-values are V-values are described by a smooth function, is also called 'function approximation in TD learning'. The family of functions can be defined by the parameters of a Neural Network or by the parameters of a linear superposition of basis functions.

# 1. Summary: Deep Neural Network for TD learning

**In all TD learning methods**

(includes n-step SARSA, Q-learning, TD( $\lambda$ ))

- V-values OR Q-values are the central quantities
- actions are taken with softmax, greedy, or epsilon-greedy policy **derived from Q-values/V-values**

(previous slide)

In the previous two weeks, we have seen many different versions of TD learning. This includes SARSA and Q-learning, TD learning, with eligibility traces (decay factor  $\lambda < 1$ ) or without, or n-step V-learning.

In all of these algorithms the V-values or Q-values are the central quantities. We first learn the V-values (or Q-values) and then the policy is based on these values.

# 1. TD learning versus Policy Gradient

**Aim for today:**

- learn actions directly
- no need for Q-value estimation

**→ Policy Gradient**

(previous slide)

The question for today is: Can we learn directly the policy – without taking the detour via the Q-values or V-values? The answer is yes and leads to a family of methods that are called ‘policy gradient’.



# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. First steps toward deep reinforcement learning
2. Basic idea of policy gradient

(previous slide)

Let us start with the reasons to work with policy gradients rather than V-values or Q-values.

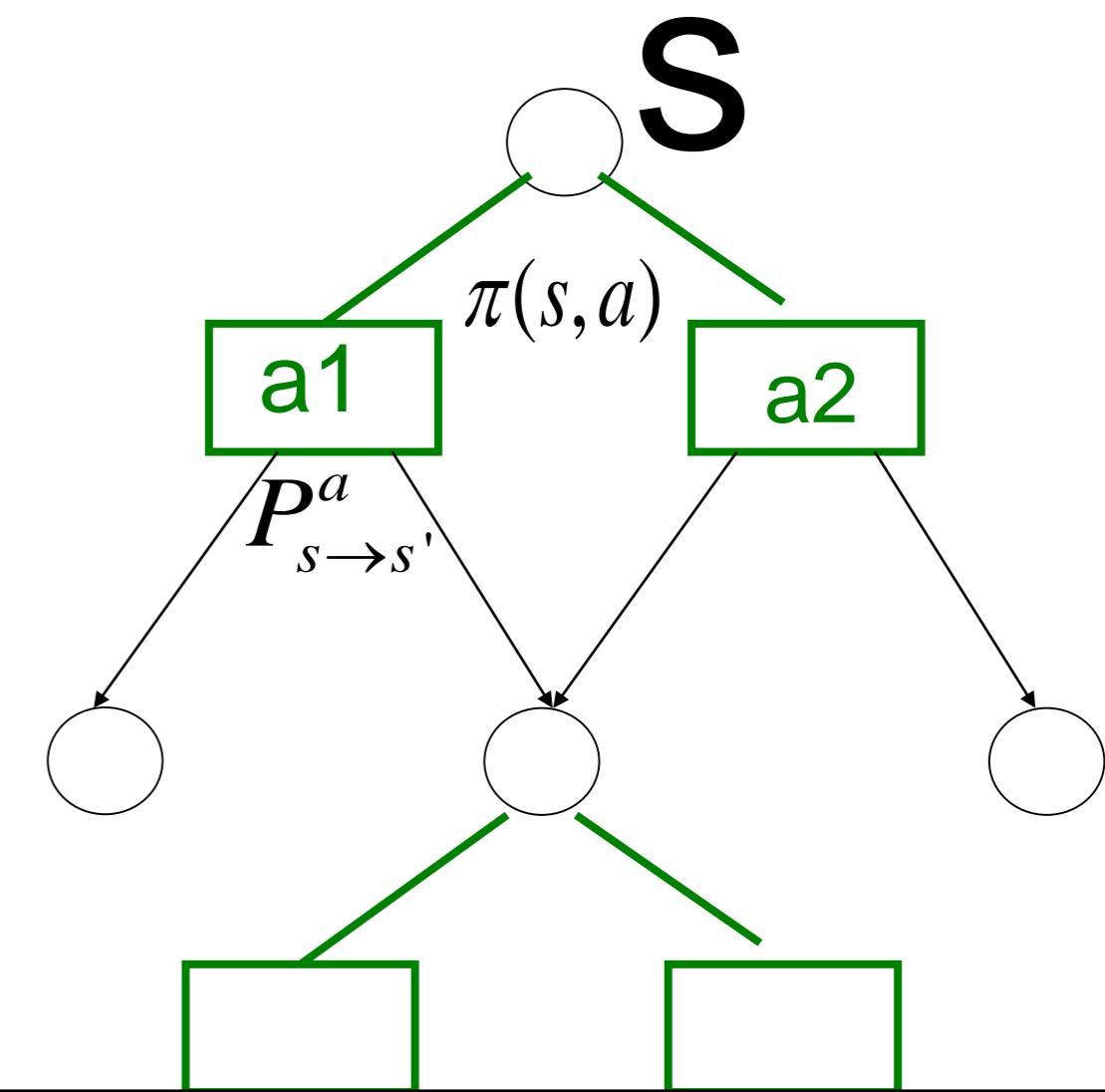
# Disadvantages of Q-learning, SARSA, or TD-learning

- For continuous states, **function approximation** is necessary (which are potentially unstable).
- Even in fully observable (Markov) settings, off-policy TD algorithms (e.g. Q-learning) can diverge using **function approximation**.
- In **partially observable** environments (non-Markov), TD algorithms are problematic
- **Continuous actions** are difficult to represent using TD methods.

*World is not a Markov Process*

*World is not fully observable*

*World is not tabular (not discrete states)*



(previous slide)

Q-values and V-values work best in an environment that has Markov properties, in particular discrete, distinguishable states, and transition probabilities between these states. Building V-values (or Q-values) then means building a table of these states (or state-action pairs).

But the world is not Markovian; however, if we use the Markovian assumptions in an environment where this is not true, then there is no guarantee that these algorithms converge.

# 1. Policy Gradient methods: basic idea

- Forget Q-values
- Optimize directly the reward
- Associate actions with stimuli stochastically

Table in Q-learning:  
(state,action)  $\rightarrow$  Q

	$a_1$	$a_2$	$a_3$
$s_1$	$Q(s_1, a_1)$	$Q(s_1, a_2)$	
$s_2$	$Q(s_2, a_1)$		
$s_3$			
$s_4$			

Table in Policy gradient:  
state  $\rightarrow$  Prob(action|state)

	$a_1$	$a_2$	$a_3$
$s_1$	0.1	0.8	0.1
$s_2$	0.75	0.1	0.15
$s_3$	0.01	0.02	0.97
$s_4$	0.5	0.5	0.0

(previous slide)

Difference between Q-learning and policy gradient:

In Q-learning you build a table of  $Q(s,a)$  for each state-action pair. Then you derive the policy from this (e.g., epsilon-greedy).

In policy gradient you learn directly the probability of taking action  $a$  in state  $s$ . Since these are probabilities, they must sum to one.

# 1. Policy Gradient methods: basic idea

- Forget Q-values
- Optimize directly the reward
- Associate actions with stimuli using a stochastic policy
- **Change parameters so as to maximize rewards**

**stochastic policy**

$$\pi(a|s, \theta)$$

parameter

Table in Policy gradient:  
state  $\rightarrow$  Prob(action|state, parameters)

	$a_1$	$a_2$	$a_3$
$s_1$	0.1	0.8	0.1
$s_2$	0.75	0.1	0.15
$s_3$	0.01	0.02	0.97
$s_4$	0.5	0.5	0.0

(previous slide)

The basic ideas are now that

(i) these probabilities will depend on a set of parameters  $\theta$

(ii) these probabilities can be directly interpreted as the policy  $\pi(a|s, \theta)$

Note sometimes the policy is written with parameters suppressed, or parameters added as an index:

$$\pi(a|s, \theta) \rightarrow \pi_{\theta}(a|s)$$



# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

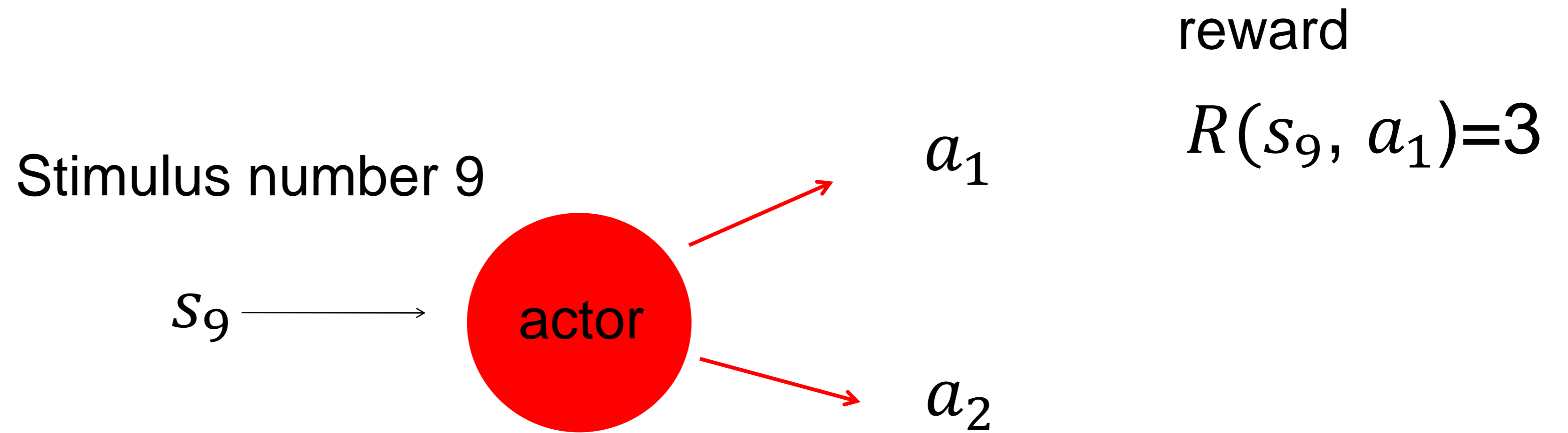
1. First steps toward Deep RL
2. Basic idea of policy gradient
3. Example: 1-step horizon

(previous slide)

To make these abstract notions concrete, we start with a simple example.

### 3. Policy Gradient methods: 1-step horizon

- Associate actions with stimuli
- Optimize directly the reward



(previous slide)

As always in reinforcement learning, the goal is to optimize rewards. We start with a one-step horizon and a binary choice.

For each stimulus (here stimulus number 9) there is the choice of two actions.

For example if the agent takes action  $a_1$  in response to stimulus  $s_9$ , it receives a reward of value 3.

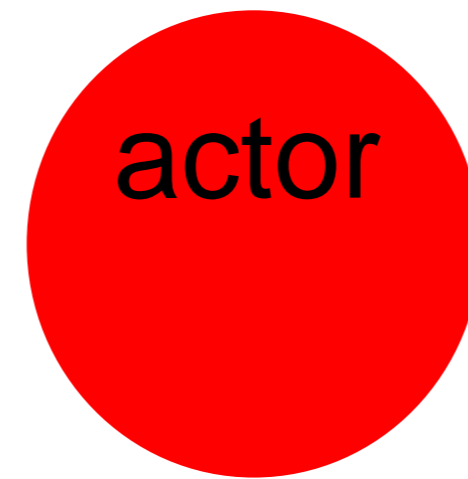
### 3. Policy Gradient methods: 1-step horizon

stimulus=state=input vector

Stimulus number 9 is a vector

$$\vec{x}^9 = (x_1^9, x_2^9, \dots, x_N^9)^T$$

$$s_9 = \vec{x}^9 \longrightarrow$$



$a_1$

$a_2$

$$R(\vec{x}^9, a_1) = 3$$

(previous slide)

We model the stimulus  $s$  as in input vector (input pattern  $\vec{x}$  ).

The actor can take two possible actions.

### 3. Policy Gradient methods: 1-step horizon

Aim: change weights of neuron

→ Maximize expected reward!

$$\langle R \rangle = \sum_x \sum_{y=\{0,1\}} \pi(y|\vec{x}) p(\vec{x}) R(y, \vec{x})$$

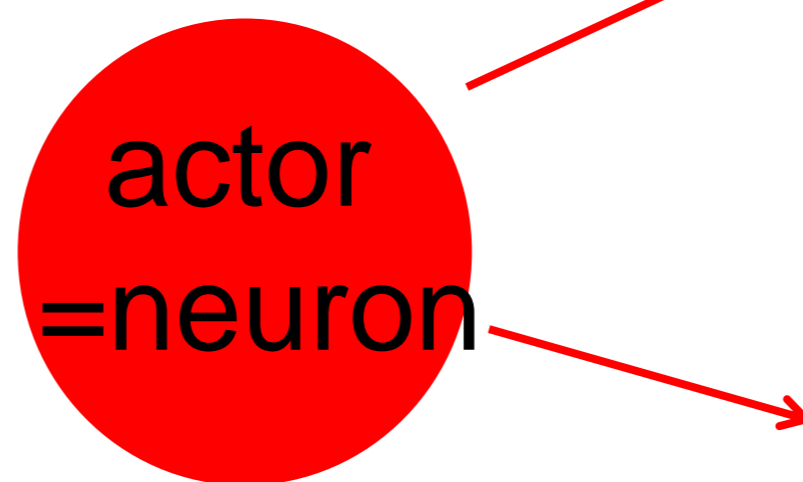
Stimulus number 9

$$\vec{x}^9 = (x_1^9, x_2^9, \dots, x_N^9)^T$$

Output of neuron

$$a_1 \rightarrow y = 1$$

$$s_9 = \vec{x}^9 \longrightarrow$$



$$a_2 \rightarrow y = 0$$

#### Choice of actions

policy:  $\pi(a_1|\vec{x}, \vec{w}) = \text{prob}(y = 1|\vec{x}, \vec{w}) = g\left(\sum_k^N w_k x_k\right)$

(previous slide)

We now model the policy as a single sigmoidal neuron with transfer function  $g$  and weight vector  $\vec{w}$ .

The question now is:

How should we adapt the weight vector so that (averaged over all possible stimuli) the reward is maximal?

Define the mean reward as

$$\langle R \rangle = \sum_{\vec{x}} \sum_{y=\{0,1\}} \pi(y|\vec{x}) p(\vec{x}) R(y, \vec{x})$$

and use  $\pi(y = 1|\vec{x}) = g(\sum_k^N w_k x_k)$



# Exercise 1: maximize expected reward

*Exercise 1  
now (10min)*

Exercise 1. (in Class): Single neuron as an actor

Assume an agent with binary actions  $Y \in \{0, 1\}$ . Action  $y = 1$  is taken with a probability  $\pi(Y = 1|\vec{x}; \vec{w}) = g(\vec{w} \cdot \vec{x})$ , where  $\vec{w}$  are a set of weights and  $\vec{x}$  is the input signal that contains the state information. The function  $g$  is monotonically increasing and limited by the bounds  $0 \leq g \leq 1$ .

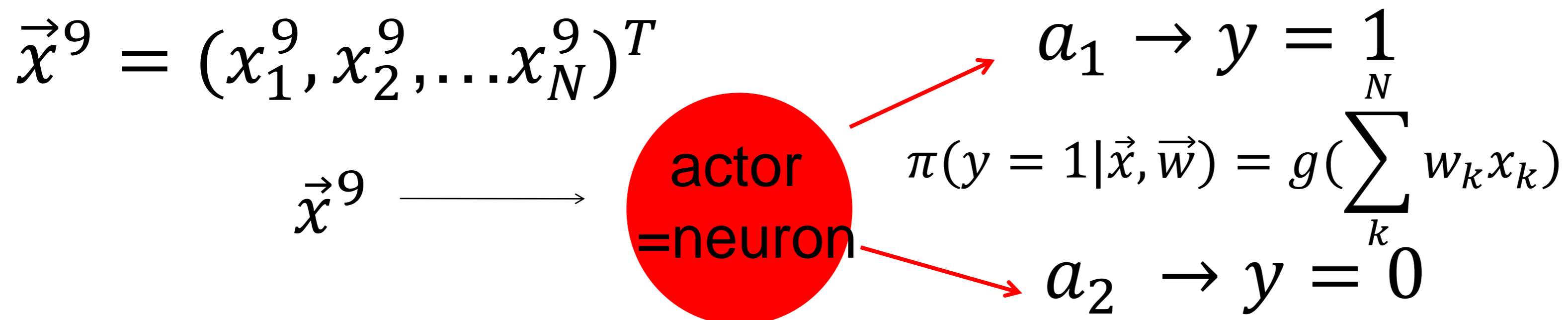
For each action, the agent receives a reward  $R(Y, \vec{x})$ .

a. Calculate the gradient of the mean reward  $\langle R \rangle = \sum_{Y, \vec{x}} R(Y, \vec{x}) \pi(Y|\vec{x}; \vec{w}) P(\vec{x})$  with respect to the weight  $w_j$ .

Hint: Insert the policy  $\pi(Y = 1|\vec{x}; \vec{w}) = g(\sum_k w_k x_k)$  and  $\pi(Y = 0|\vec{x}; \vec{w}) = 1 - g(\sum_k w_k x_k)$ . Then take the gradient.

b. The rule derived in (a) is a batch rule. Can you transform this into an 'online rule'?

Hint: Pay attention to the following question: what is the condition that we can simply 'drop the summation signs'?



(your calculations)

# 3. Policy Gradient methods: 1-step horizon

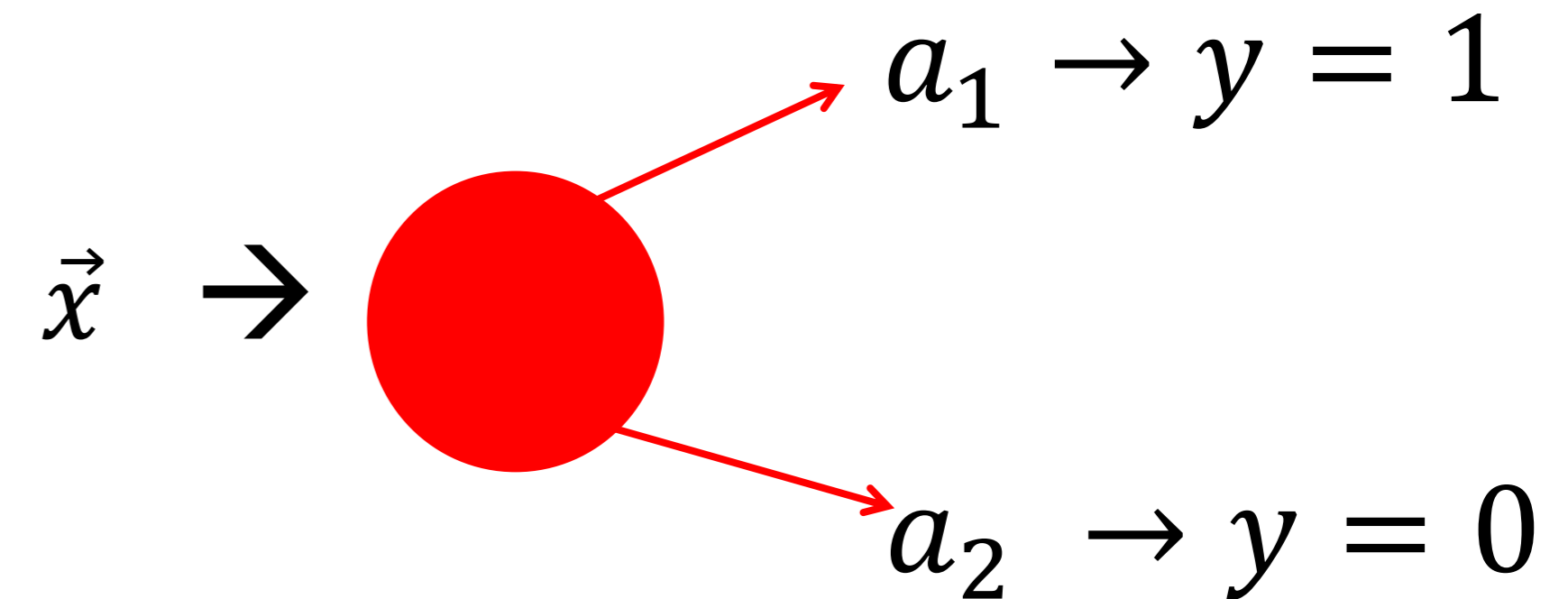
*blackboard1*

reward

$$R(y, \vec{x})$$

policy

$$\pi(y = 1 | s, \vec{w}) = g\left(\sum_k^N w_k x_k\right)$$



(previous slide)

It is convenient to introduce a binary output variable:  $y$  takes the value of 1 if action  $a_1$  is taken and zero otherwise.

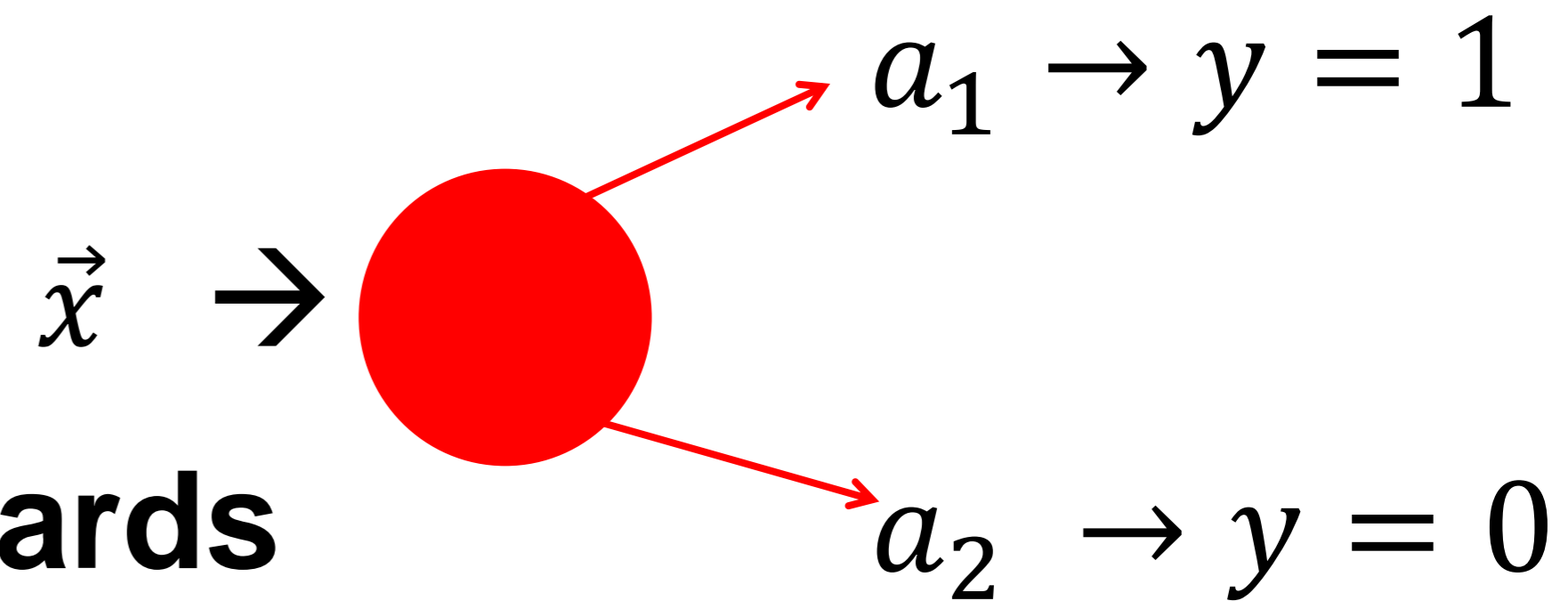
# 3. Policy Gradient methods: 1-step horizon

**reward**

$$R(y, \vec{x})$$

**policy**

$$\pi(y = 1 | s, \vec{w}) = g\left(\sum_k^N w_k x_k\right)$$



**Update parameters to maximize rewards**

$$\text{If } y = 1: \quad \Delta w_j = \eta \frac{g'}{g} R(1, \vec{x}) x_j$$

$$\text{If } y = 0: \quad \Delta w_j = \eta \frac{-g'}{(1-g)} R(0, \vec{x}) x_j$$

(previous slide)

The optimal update rule (last two lines) has a simple interpretation:

The weight  $w_j$  is moved in direction of  $x_j$  if the reward is positive.

The notation  $g'$  refers to the derivative of the sigmoidal function  $g$ .

# 3. Policy Gradient methods: Batch-to-Online

Attention at transition 'Batch to Online':  
→ natural statistical weight must be correct!

We have a stochastic starting point with weight  $p(s)$   
as well as stochastic transitions and a stochastic policy

$$\sum_{s'} P_{s \rightarrow s'}^a$$

weighting factor  
for 'next state'

$$\sum_{a'} \pi(a' | s, \vec{w})$$

weighting factor  
for 'next action'

(previous slide)

Batch rule (like in standard ANN): a single update is performed after having processed many patterns (minibatch) or all patterns (standard batch rule).

Online rule (like in standard ANN): an update is performed at every time step (after each pattern).

The example (and your calculations in the exercise) show that the transition from batch to online is not always possible by deleting the sum signs. In fact, it is only possible if the statistical weighting factor is correct.



# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. First steps toward Deep RL
2. Basic idea of policy gradient
3. Example: 1-step horizon
4. Log-likelihood trick

(previous slide)

Is there a more systematic way to perform the transition from batch to online?

The answer is yes and given by (what I call) the log-likelihood trick

# 4. Log-likelihood trick

*Blackboard 2*

(your comments)

## 4. Log-likelihood trick

$$\begin{aligned}\nabla_{\theta} J &= \int \nabla_{\theta} p(H) R(H) dH \\ &= \int \frac{p(H)}{p(H)} \nabla_{\theta} p(H) R(H) dH \\ &= \int p(H) \nabla_{\theta} \log p(H) R(H) dH\end{aligned}$$

J = function you want to optimize

H = ensemble over which you integrate

(previous slide)

From BATCH to ONLINE

Suppose you want to optimize some function  $J$  which is given by the integral over the statistical ensemble  $H$ . Instead of an integral you often have the sum over all possible patterns, for example.

You want to do optimization by gradient ascent, therefore you need to calculate the gradient.

For the correct statistical weight you need the weight factor  $p(H)$ .

Normally this factor disappears when you naively take the gradient. However, if you rewrite this as the gradient of  $(\log p)$  and then multiply by  $p(H)$ , you have the exactly the same result – but now the correct weight factor  $p(H)$  is explicit. Now you can cut out the integral and  $p(H)$  and you get a valid online rule.

## 4. Policy gradient derivation

$$\nabla_{\theta} J = \int p(H) \nabla_{\theta} \log p(H) R(H) dH.$$

Taking the sample average as Monte Carlo (MC) approximation of this expectation by taking  $N$  trial histories we get

$$\nabla_{\theta} J = \mathbf{E}_H \left[ \nabla_{\theta} \log p(H) R(H) \right] \approx \frac{1}{N} \sum_{n=1}^N \nabla_{\theta} \log p(H^n) R(H^n).$$

which is a fast approximation of the policy gradient for the current policy

(previous slide)

Delete the integral and  $p(H)$  and sum over all examples, and you have a good approximation to your original integral.



# 4. Policy gradient evaluation: Example from Exercise 1

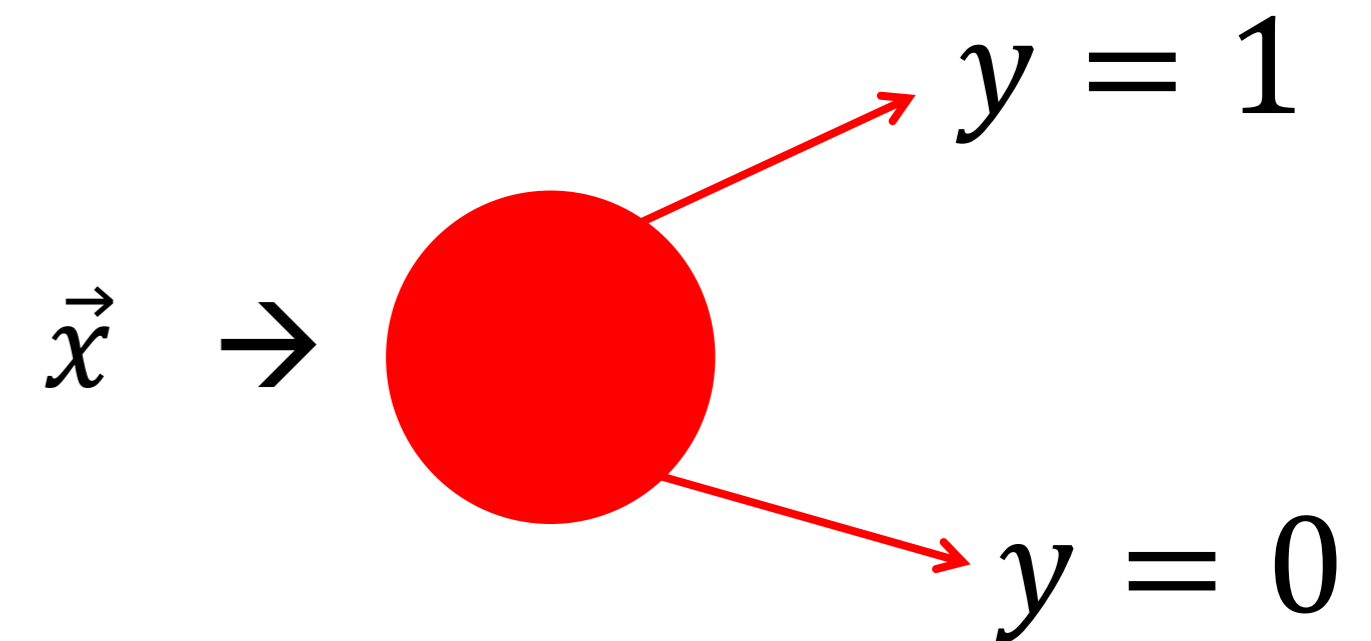
reward

*Blackboard 2b*

$$R(y, \vec{x})$$

policy

$$\pi(y = 1 | s, \vec{w}) = g\left(\sum_k^N w_k x_k\right)$$



(your comments)

# 4. Update rule for Exercise 1

observe input  $\vec{x}$ , output  $y$ , and reward  $R(y, \vec{x})$

**Earlier result:**

$$\text{If } y = 1: \quad \Delta w_j = \eta \frac{g'}{g} \cdot R(1, \vec{x}) x_j$$

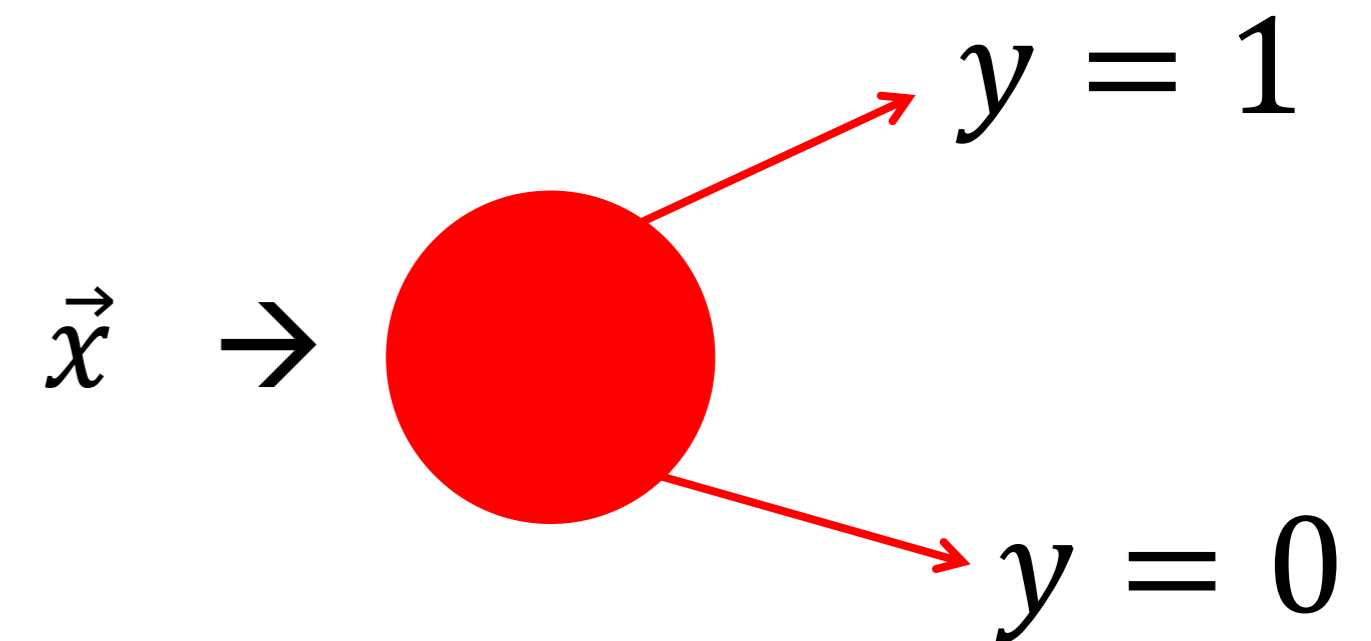
$$\text{If } y = 0: \quad \Delta w_j = \eta \frac{-g'}{(1-g)} R(0, \vec{x}) x_j$$

**policy**

$$\pi(y = 1 | s, \vec{w}) = g \left( \sum_k^N w_k x_k \right)$$

**Now rewritten as:**

$$\Delta w_j = \eta \frac{g'}{g(1-g)} R(y, \vec{x}) [y - \langle y \rangle] x_j$$



Note:  $\langle y \rangle = g(\sum_k^N w_k x_k)$

(previous slide)

Using the log-likelihood trick we arrive at the same result as before but faster and, importantly, via a systematic sequence of steps.

Last line – two important comments:

- (i) The two cases ( $y=+1$ ) and ( $y=0$ ) can be summarized in a single update rule
- (ii)  $\langle y \rangle$  is the expectation of the output, given the input vector  $\vec{x}$

## 4. Comparison with Perceptron

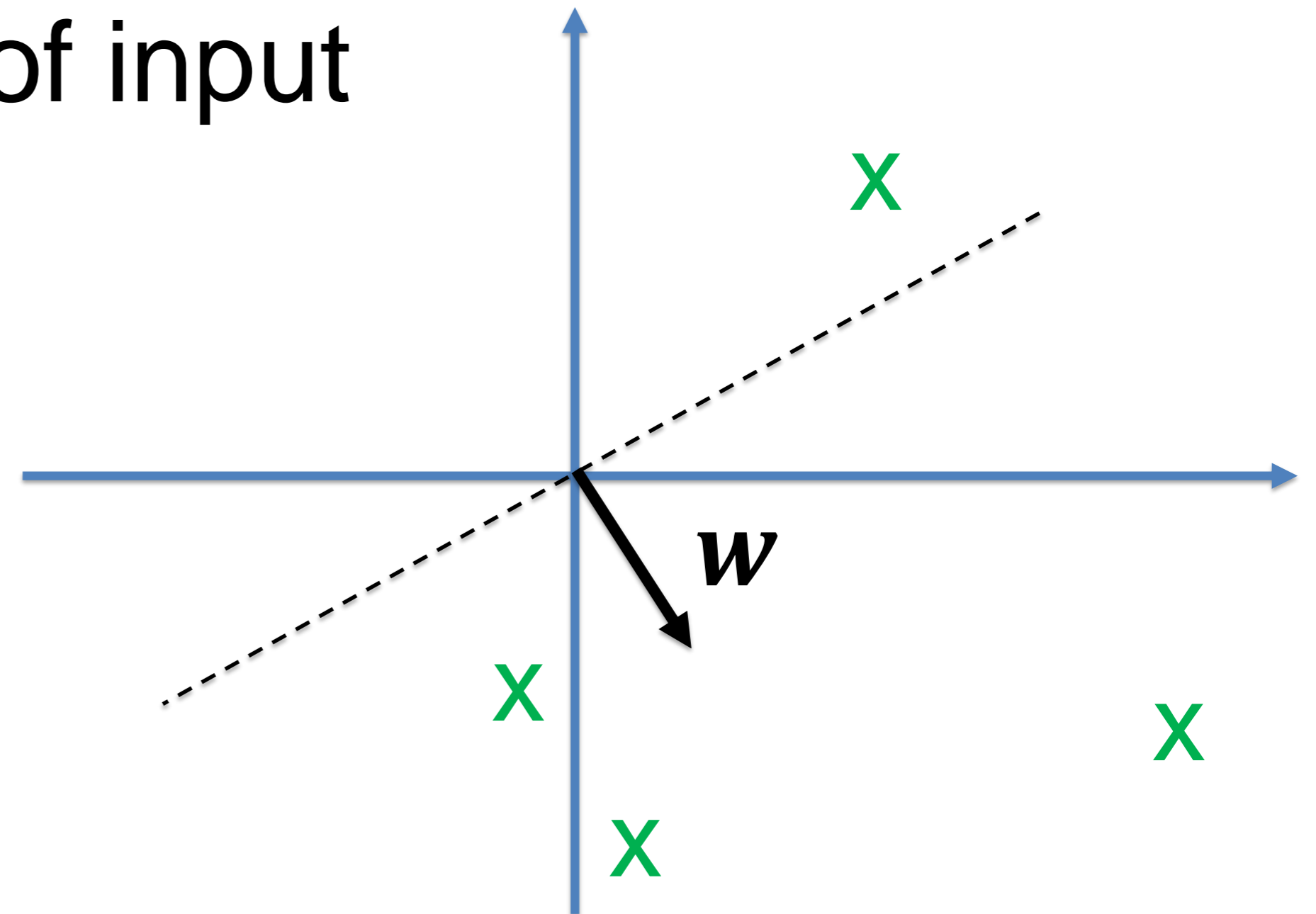
*parameter = weight  $w_j$*

$$\Delta w_j \propto R(y, \vec{x}) [y - \langle y \rangle] x_j$$

Weight vector turns in direction of input

$$\Delta \mathbf{w} \propto \pm \mathbf{x}$$

$$R > 0 \text{ and } y = 1 \rightarrow \Delta \mathbf{w} \propto +\mathbf{x}$$



(previous slide)

Similar to the perceptron update rule, the update with gradient descent can be interpreted as a weight vector that turns in direction of an input pattern (with positive or negative sign)

## 4. Comparison with Biology

*parameter = weight*  $w_j$

$$\Delta w_j \propto R(y, \vec{x}) [y - \langle y \rangle] x_j$$

Stimulus

reward

pre

j

i

post

Weight vector turns in direction of input

Three factors: reward

post

pre

$$\Delta w_{ij} = \eta \frac{g'}{g(1-g)} R(\vec{y}, \vec{x}) [y_i - \langle y_i \rangle] x_j$$

postsynaptic factor is

*'activity – expected activity'*

(previous slide)

The update rule give also rise to an interesting biological interpretation.

The learning rule depends on three factors:

- (i) The reward
- (ii) The 'state' of the postsynaptic neuron where 'state'=activity minus expected activity
- (iii) The presynaptic activity

Presynaptic cell: the neuron that sends a signal across the connection (sender)

Postsynaptic cell: the neuron that receives the signal (receiver).

We will come back to the link between reinforcement learning and the brain in lecture 12.



## 4. Generalization: subtract a reward baseline

we derived this online gradient rule

$$\Delta w_j \propto R(y, \vec{x}) [y - \langle y \rangle] x_j$$

But then this rule is also an online gradient rule

$$\Delta w_j \propto [R(y, \vec{x}) - b] [y - \langle y \rangle] x_j$$

with the same expectation (see exercise 2)

(because a baseline shift drops out if we take the gradient)

(previous slide)

Note that we are interested in finding the set of weights that optimize the expected reward  $\langle R \rangle$ .

The update rule has been derived by taking the gradient on the mean reward  $\langle R \rangle$ .

But a function  $\langle R - b \rangle$  with constant bias  $b$  would have exactly the same location of the maximum.

If we repeated the gradient steps, the results would lead to an update rule with a factor  $[R - b]$  instead of  $R$ . Therefore, the rule with  $[R - b]$  is also a valid online rule.

## 4. Quiz: Policy Gradient and Reinforcement learning

Your friend has followed over the weekend a tutorial in reinforcement learning and claims the following. Is he right?

- All reinforcement learning algorithms work either with Q-values or V-values
- The transition from batch to online is always easy: you just drop the summation signs and bingo!
- All reinforcement learning algorithms try to optimize the expected total reward (potentially discounted if there are multiple time steps)
- The derivative of the log-policy is some abstract quantity that has no intuitive meaning.

(your comments)

# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. First steps toward Deep RL
2. Basic idea of policy gradient
3. Example: 1-step horizon
4. Log-likelihood trick
5. Multiple time steps

(previous slide)

So far the discussion has been restricted to scenarios with a one-step horizon.

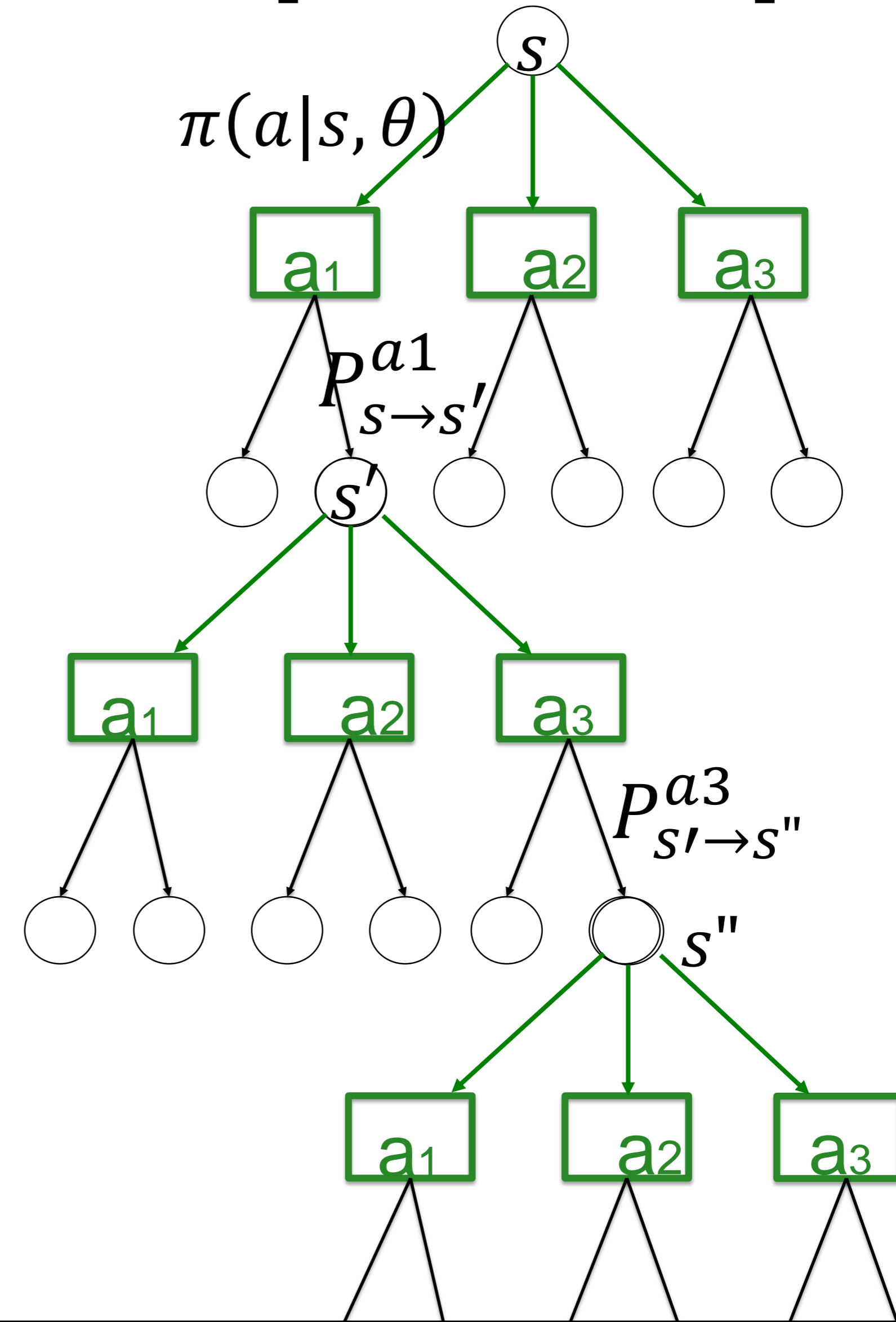
The agent takes an action, gets a reward, and the episode ends.

Now we need to generalize to scenarios that extend over multiple time steps.

# 5. Policy Gradient methods over multiple time steps

Aim:

update the parameters  $\theta$   
of the policy  $\pi(a|s, \theta)$



(previous slide)

We use the same graph of the multistep Markov decision model as for the derivation of the Bellman equation.

However, now we work directly on a policy  $\pi(a|s,\theta)$  which depends on parameters  $\theta$ .



# 5. Policy Gradient methods over multiple time steps

*Blackboard 3*

(your notes)

# 5. Policy Gradient methods over multiple time steps

Calculation yields several terms of the form

Total accumulated discounted reward  
collected in one episode starting at  $s_t, a_t$

$$\Delta\theta_j \propto \left[ R_{s_t \rightarrow s_{end}}^{a_t} \right] \frac{d}{d\theta_j} \ln[\pi(a_t | s_t, \theta)]$$
$$+ \gamma \left[ R_{s_{t+1} \rightarrow s_{end}}^{a_{t+1}} \right] \frac{d}{d\theta_j} \ln[\pi(a_{t+1} | s_{t+1}, \theta)]$$
$$+ \dots$$

(previous slide)

We consider a single episode that started in state  $s_t$  with action  $a_t$  and ends after several steps in the terminal state  $s_{end}$

The result of the calculation gives an update rule for each of the parameters.

The update of the parameter  $\theta_j$  contains several terms.

(i) the first term is proportional to the total accumulated (discounted) reward, also called return  $R_{s_t \rightarrow s_{end}}^{a_t}$

(ii) the second term is proportional to gamma times the total accumulated (discounted) reward but starting in state  $s_{t+1}$

(iii) the third term is proportional to gamma-squared times the total accumulated (discounted) reward but starting in state  $s_{t+2}$

(iv)

We can think of this update as one update step for one episode. Analogous to the terminology last week, Sutton and Barto call this the Monte-Carlo update for one episode.

Note that each of the terms is proportional to  $\ln \pi$

# Policy Gradient methods over multiple time steps:

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for  $\pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Algorithm parameter: step size  $\alpha > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (G_t)$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \theta)$$

From book:

Sutton and Barto, 2018

Different states  $S_0, S_1, S_2, \dots$  during one episode

$G$  = total accumulated reward during the episode starting at  $S_t$ ;

All updates done AT THE END of the episode

Algorithm maximizes expected discounted rewards starting at  $S_0$

(previous slide)

The algorithm in Pseudocode taken from the book of Sutton and Barto. The update concerns a single episode.

The only notational difference with respect to the earlier slide is a rewrite of the factors gamma – you can check the equivalence by taking a piece of paper.

Note that for an implementation it would be most convenient to start at the terminal state of the episode and work backwards so as to reuse the return calculations.

Variations of this algorithm are the basis of policy gradient methods and widely used in applications.

# Artificial Neural Networks: RL3

## Policy Gradient Methods

Wulfram Gerstner

EPFL, Lausanne, Switzerland

1. Review
2. Basic idea of policy gradient
3. Example: 1-step horizon
4. Log-likelihood trick
5. Multiple time steps
6. Subtracting the mean via the value function

(previous slide)

In the simple one-step scenario we have seen that we can subtract a bias  $b$  from the reward.



## 6. Review: subtract a reward baseline

we derived this online gradient rule (for 1-step horizon)

$$\Delta w_j \propto R(y, \vec{x}) [y - \langle y \rangle] x_j$$

But then this rule is also an online gradient rule

$$\Delta w_j \propto [R(y, \vec{x}) - b] [y - \langle y \rangle] x_j$$

with the same expectation (see exercise 2)

(because a baseline shift drops out if we take the gradient)

(previous slide)

The question arises whether the same is true in the multi-step episodes.

The answer is YES.

## 6. Subtract a reward baseline

we derived this online gradient rule for multi-step horizon

$$\Delta\theta_j \propto [R_{s_t \rightarrow s_{end}}^{a_t}] \frac{d}{d\theta_j} \ln[\pi(a_t | s_t, \theta)]$$

But then this rule is also an online gradient rule

$$\Delta\theta_j \propto [R_{s_t \rightarrow s_{end}}^{a_t} - b(s_t)] \frac{d}{d\theta_j} \ln[\pi(a_t | s_t, \theta)]$$

with the same expectation

(because a baseline shift drops out if we take the gradient)

(previous slide)

Please remember that the full update rule for the parameter  $\theta_j$

in a multi-step episode contains several terms of this form; here only the first of these terms is shown.

Similar to the case of the one-step horizon, we can subtract a bias  $b$  from the reward without changing the location of the maximum of the total expected return.

Moreover, this bias  $b(s_t)$  can itself depend on the state  $s_t$ .

Thus the update rule now has terms

$$\begin{aligned} \Delta\theta_j \propto & [R_{s_t \rightarrow s_{end}}^{a_t} - b(s_t)] \frac{d}{d\theta_j} \ln[\pi(a_t | s_t, \theta)] \\ & + \gamma [R_{s_{t+1} \rightarrow s_{end}}^{a_{t+1}} - b(s_{t+1})] \frac{d}{d\theta_j} \ln[\pi(a_{t+1} | s_{t+1}, \theta)] \\ & + \gamma^2 [R_{s_{t+2} \rightarrow s_{end}}^{a_{t+2}} - b(s_{t+2})] \frac{d}{d\theta_j} \ln[\pi(a_{t+2} | s_{t+2}, \theta)] \\ & + \dots \end{aligned}$$

## 6. subtract a reward baseline

Total accumulated discounted reward  
collected in one episode starting at  $s_t, a_t$

$$\Delta\theta_j \propto [R_{s_t \rightarrow s_{end}}^{a_t} - b(s_t)] \frac{d}{d\theta_j} \ln[\pi(a_t | s_t, \theta)] + \dots$$

- The bias  $b$  can depend on state  $s$
- Good choice is  $b = \text{'mean of } [R_{s_t \rightarrow s_{end}}^{a_t}] \text{'}$ 
  - take  $b(s_t) = V(s_t)$
  - learn value function  $V(s)$

(previous slide

Is there a choice of the bias  $b(s_t)$  that is particularly good?

One attractive choice is to take the bias equal to the expectation (or empirical mean). The logic is that if you take an action that gives more accumulated discounted reward than your empirical mean in the past, then this action was good and should be reinforced.

If you take an action that gives less accumulated discounted reward than your empirical mean in the past, then this action was not good and should be weakened.

But what is the expected discounted accumulated reward? This is, by definition, exactly the value of the state. Hence a good choice is to subtract the V-value.

And here is where finally the idea of Bellman equation and TD learning comes in through the backdoor: we can learn the V-value, and then use it as a bias in policy gradient.

# 6. Deep reinforcement learning: alpha-zero

Network for choosing action

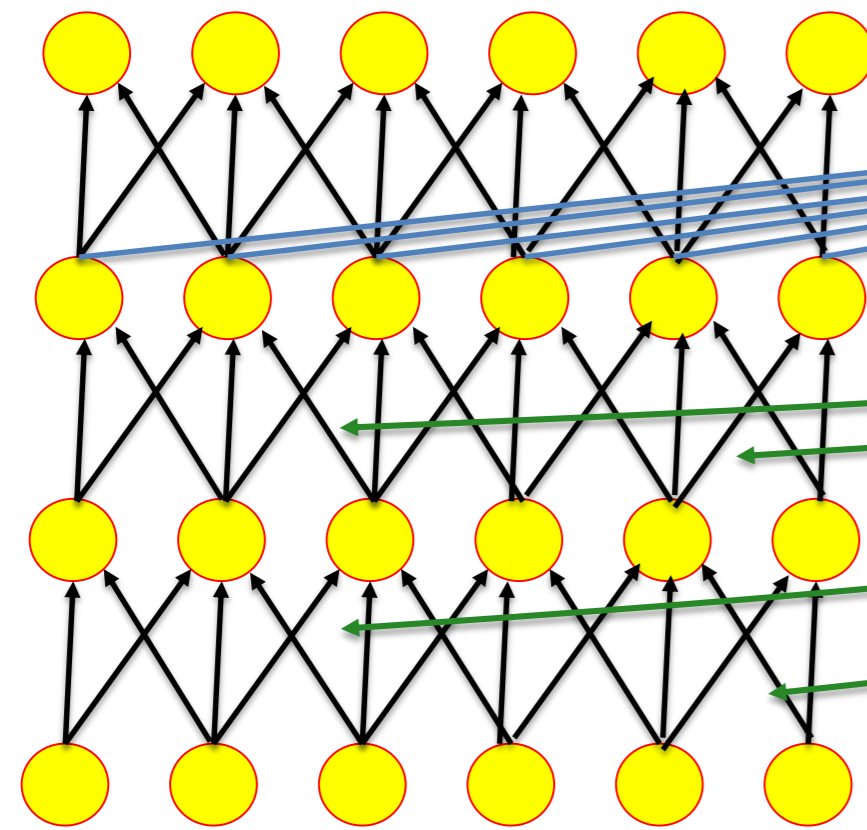
action:

*Advance king*

2<sup>e</sup> output for **value** of state:

output ↑ ↑ ↑ ↑ ↑

$V(s)$



**learning:**

→ change connections

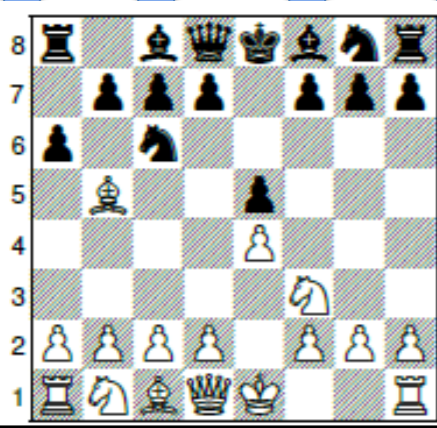
**aims:**

- learn value  $V(s)$  of position
- learn action policy to win

**Learning signal:**

- $\eta[\text{actual Return} - V(s)]$

input



(previous slide)

Very schematically is this one of the ideas of deep reinforcement learning. We construct a deep network with multiple layers. We use the output units for action choice and optimize the parameters via policy gradient. We have a further output unit to estimate the  $V$ -value, and use it as a bias.

The model of the  $V$ -value can share some units with the model of the actions ...



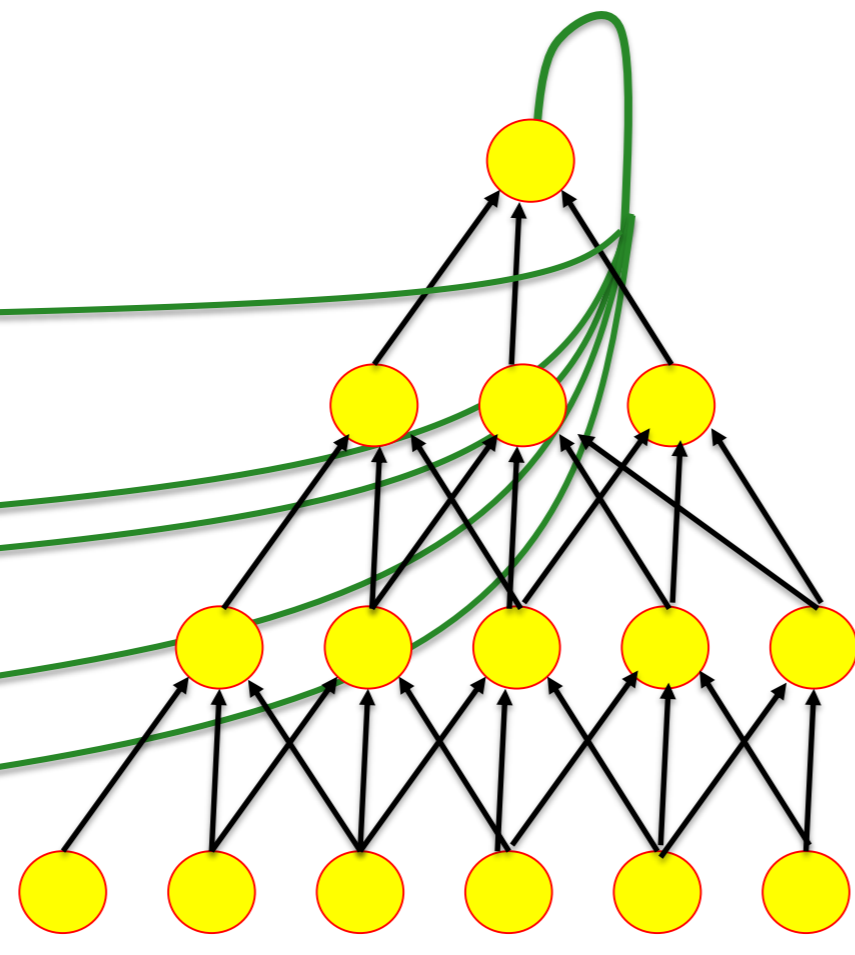
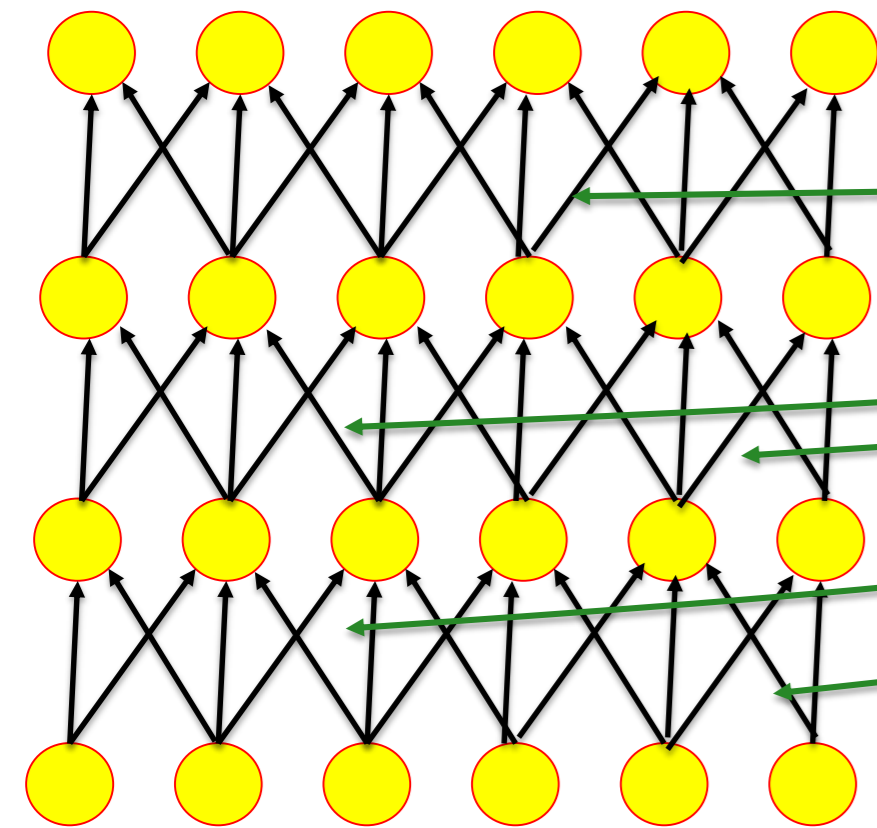
# 6. Deep Reinforcement Learning: Lunar Lander and other games (miniproject)

actions

*advance*      *push left*

value

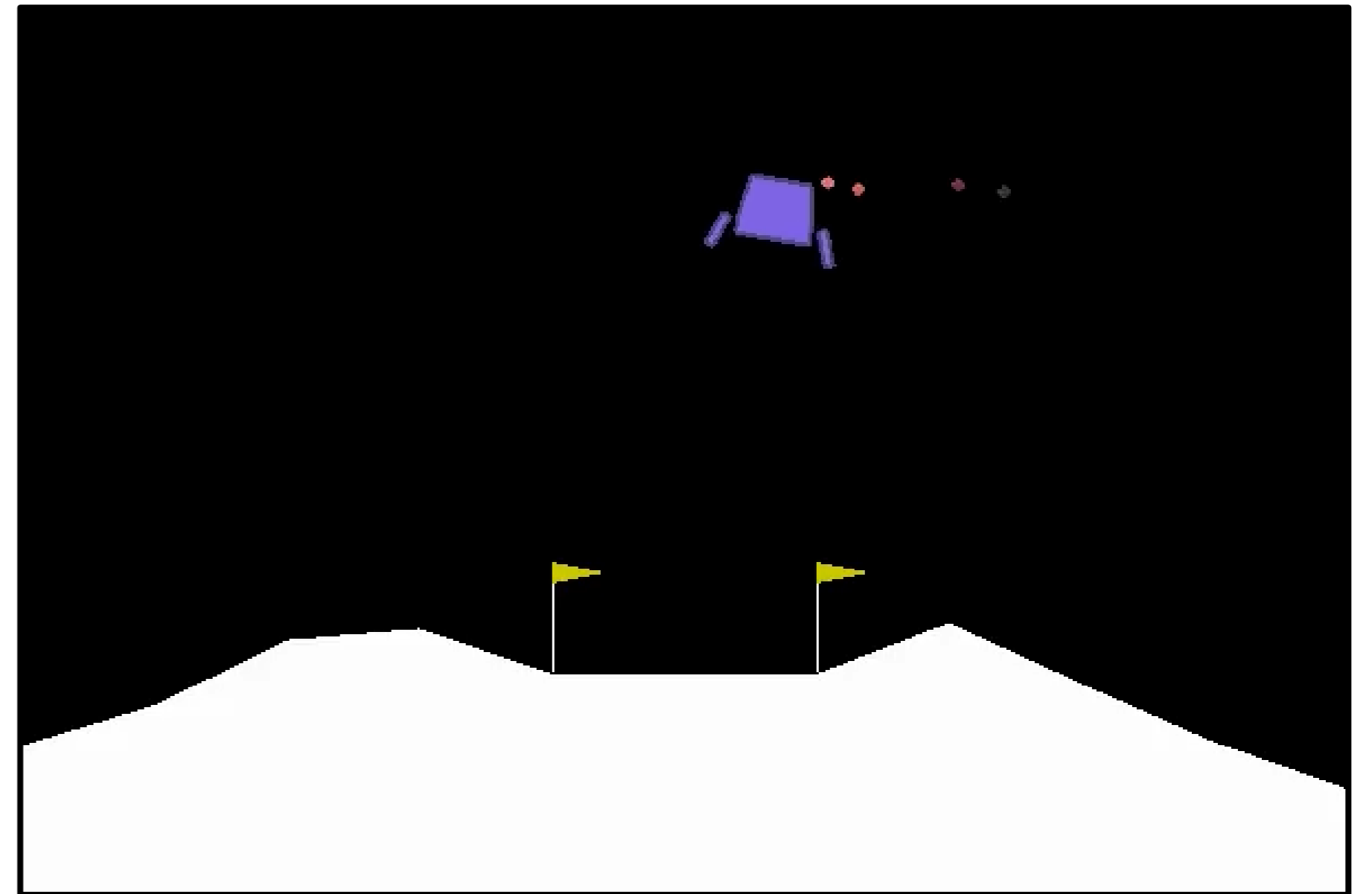
$V(s)$



$x$

$x$

Aim: land between poles



(previous slide)

... or it can be estimated in a separate network.

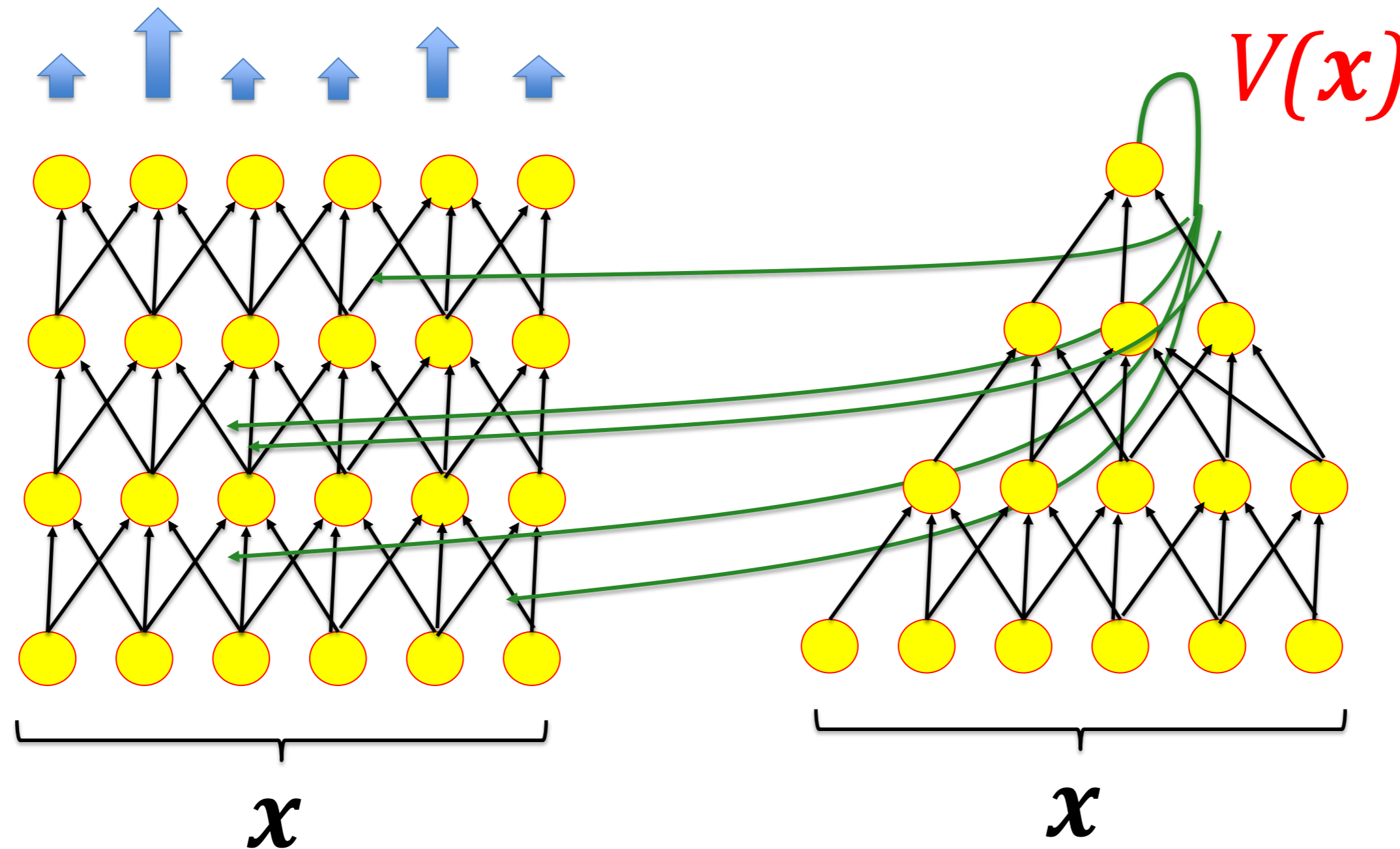
# 6. Learning two Neural Networks: actor and value

## Actions:

- Learned by Policy gradient
- Uses  $V(x)$  as baseline

## Value function:

- Estimated by Monte-Carlo
- provides baseline  $b=V(x)$  for action learning



$x$  = states from episode:

$S_t, S_{t+1}, S_{t+2},$

(previous slide)

In the latter case we have two networks:

The actor network learns a first set of parameters, called  $\theta$  in the algorithm of Sutton and Barto.

The value network learns a second set of parameters, with the label  $w$ .

The value  $b(x = s_{t+n}) = V(x)$  is the estimated total accumulated discounted reward of an episode starting at  $x = s_{t+n}$

The total accumulated discounted ACTUAL reward in ONE episode is  $R_{s_{t+n} \rightarrow s_{end}}^{a_{t+n}}$

# 'REINFORCE' with baseline

From book:  
Sutton and Barto, 2018

REINFORCE with Baseline (episodic), for estimating  $\pi_{\theta} \approx \pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$

Algorithm parameters: step sizes  $\alpha^{\theta} > 0, \alpha^{\mathbf{w}} > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \quad (G_t)$$

$$\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \gamma^t \delta \nabla \hat{v}(S_t, \mathbf{w})$$

$$\theta \leftarrow \theta + \alpha^{\theta} \gamma^t \delta \nabla \ln \pi(A_t | S_t, \theta)$$

(previous slide)

Algorithm in pseudocode taken from the book of Sutton and Barto.

For the actor, the algorithm evaluates terms of the form

$$\left[ R_{s_{t+n} \rightarrow s_{end}}^{a_{t+n}} - b(s_{t+n}) \right] \frac{d}{d\theta_j} \ln[\pi(a_{t+n} | s_{t+n}, \theta)]$$

Where the return is  $G = R_{s_{t+n} \rightarrow s_{end}}^{a_{t+n}}$

And the bias estimate is  $v(s_{t+n}) = b(s_{t+n})$

The terminal state in their notation occurs at time T and the initial state has index 0.

For the value function, they use Monte-Carlo estimation of the total accumulated reward in one episode (see last week).

## 6. Why subtract the mean?

Subtracting the expectation provides estimates that have (normally) smaller variance (look less noisy)  
→ exercise 2.

(previous slide)

Why is it useful to subtract the mean?

Whatever the choice of bias, the algorithm should eventually converge to the same set of parameters.

However, since the algorithm is based on stochastic gradient descent (i.e., the online rule instead of the full batch rule), the algorithm makes noisy steps that only go on average in the right direction.

Subtracting a bias that is close to the mean generally reduces the noise.

(Unfortunately, the minimal noise is not exactly the situation where one subtracts the mean, but it is close to it).

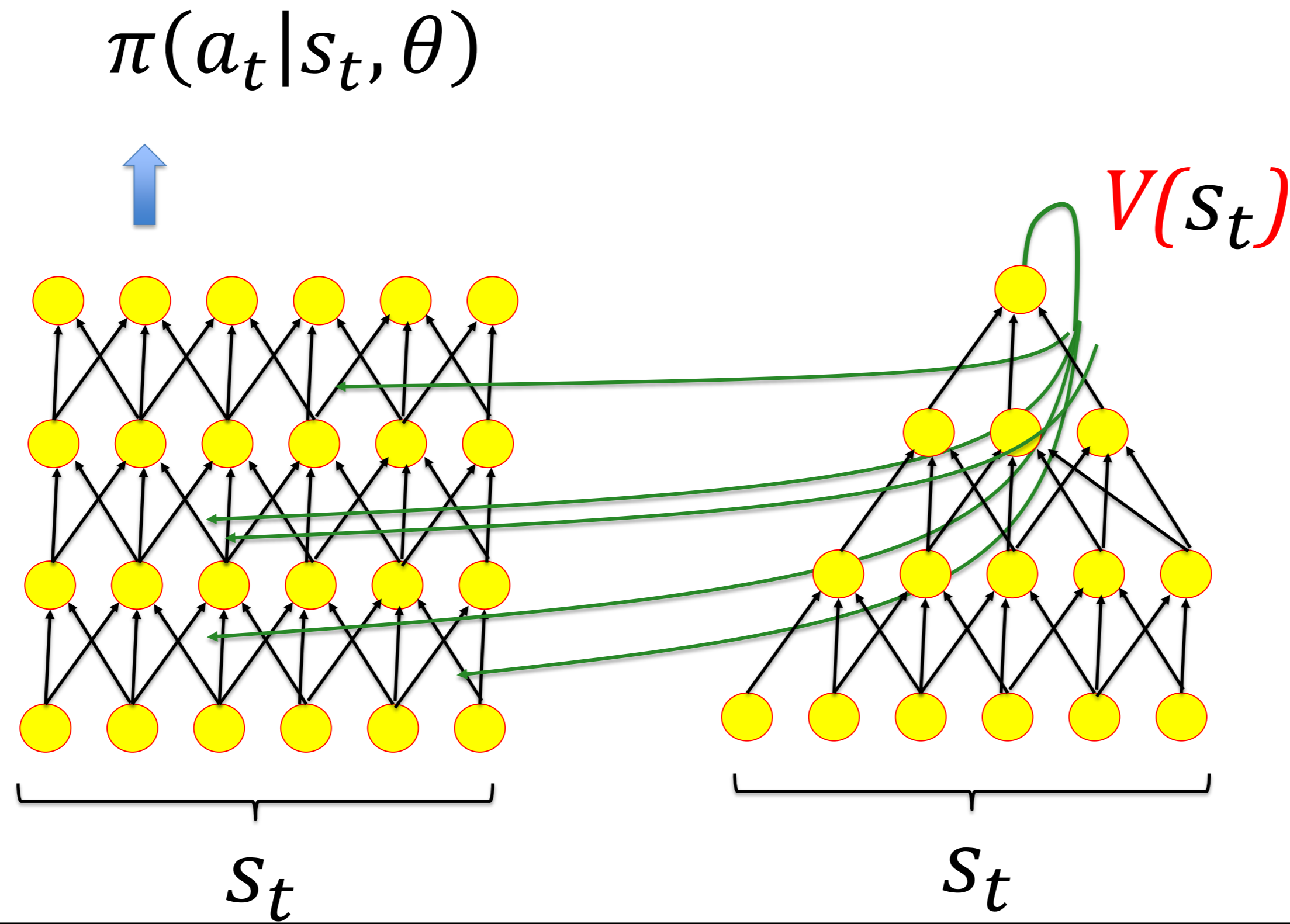


Policy gradient involves a term:

$$\frac{d}{d\theta_j} \ln[\pi(a_t | s_t, \theta)]$$

Parameters are the weights of the network.

NEXT WEEK: BackProp to calculate gradient in deep networks



(previous slide)

Both actor and critic are optimized by changing the parameters according to a gradient descent rule.

Gradient descent in a multi-layer network is called BackProp or Deep Learning.

We start with Deep Learning next week. Algorithm in pseudocode taken from the book of Sutton and Barto.

For the actor, the algorithm evaluates terms of the form

## 4. Quiz: Policy Gradient and Reinforcement learning

Your friend has followed over the weekend a tutorial in reinforcement learning and claims the following. Is he right?

Even some policy gradient algorithms use V-values

V-values for policy gradient are calculated in a separate network (but some parameters can be shared with the actor network)

(previous slide) Your notes

# Learning outcomes and Conclusions:

- **basic idea of policy gradient: learn actions, not Q-values**
  - gradient ascent of total expected discounted reward
- **log-likelihood trick: getting the correct statistical weight**
  - enables transition from batch to online
- **policy gradient algorithms**
  - updates of parameter propto  $[R - V] \frac{d}{d\theta_j} \ln[\pi]$
- **why subtract the mean reward?**
  - reduces noise of the online stochastic gradient
- **Reinforce with baseline**
  - a further output to subtract the mean reward

$$[R(s) - V(s)] \frac{d}{d\theta_j} \ln[\pi]$$

(your comments)

**The END**