

# Artificial Neural Networks: RL1 (continued)

Wulfram Gerstner  
EPFL, Lausanne, Switzerland

## Reinforcement Learning and SARSA

Parts 4-6: Examples of Reward-based Learning

### Objectives for Lecture RL1:

- Reinforcement Learning (RL) is learning by rewards ✓
- Agents and actions, states and rewards ✓
- Convergence in expectation
- Exploration vs Exploitation
- Bellman equation
- SARSA algorithm

Reading:

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

Chapters: 1.1-1.4; 2.1-2.6; 3.1-3.5; 6.4

## Reading for this week:

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

Chapters: 1.1-1.4; 2.1-2.6; 3.1-3.5; 6.4

## Background reading:

Silver et al. 2017, Archive

*Mastering Chess and Shogi by Self-Play with a  
General Reinforcement Learning Algorithm*

# Recall: Learning by reward



Learning by reward “  
BUT:  
Reward is rare:  
‘sparse feedback’ after  
a long action sequence



Previous slide.

How does a human learn to play table tennis: How does a child learn to play the piano? How does a dog learn to perform tricks?

In all these cases there is no supervisor. No master guides the hand of the players during the learning phase. Rather the player 'discovers' good movements by rather coarse feedback. For example, the ball in table tennis does not land on the table as it should. That is bad (negative feedback). The ball has a great spin so that the opponent does not get. This is good (positive feedback).

Similarly, it is hard to tell a dog what to do. But if you reinforce the dog's behavior by giving a 'goodie' at the moment when it spontaneously performs a nice action, then it can learn quite amazing things.

In all these cases it is the 'reward' that guides the learning. Rewards can be the goodie for the dog, or just the feeling 'now I did well' for humans.

# Recall: Value is important

Network for choosing action

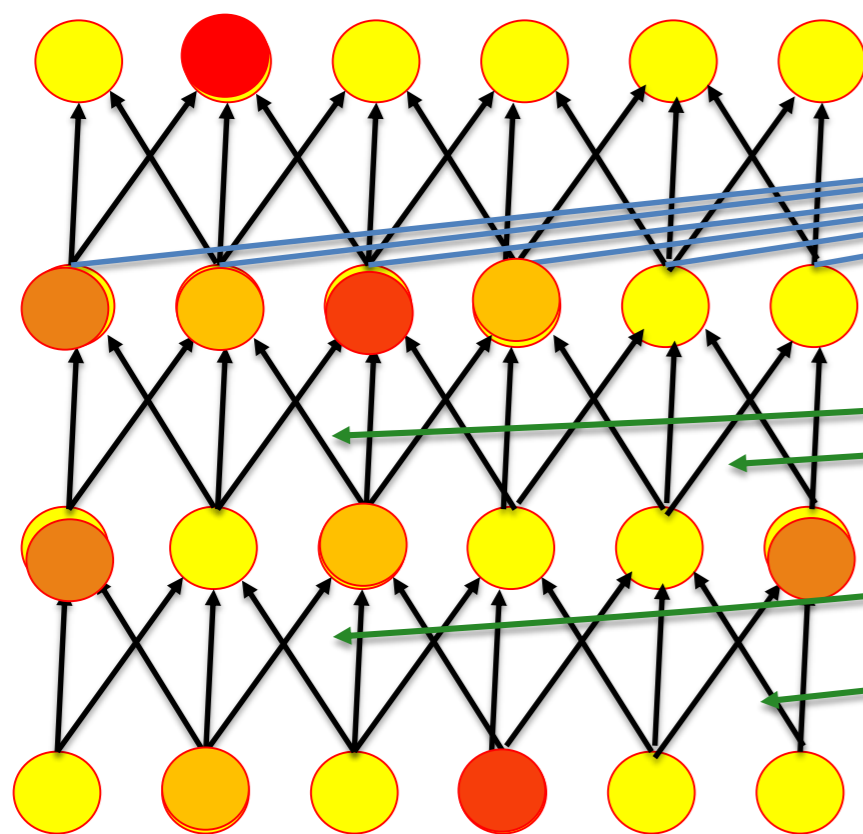
action:

*Advance king*

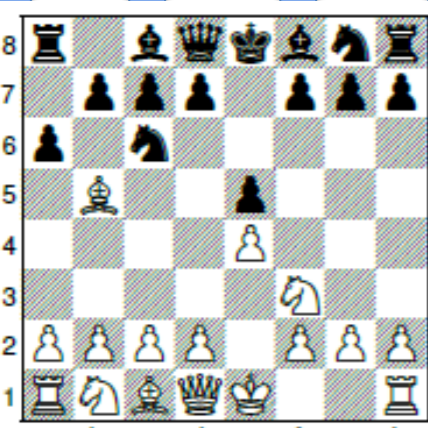
2<sup>nd</sup> output for **value** of state:

*probability to win = expected reward*

output ↑ ↑ ↑ ↑ ↑



input



**Learning by success signal**

- change connections

**aim:**

- choose next action to win

**aim for value unit:**

- predict value of current

position

Previous slide.

At the end of this semester, you will be able to understand the algorithms and network structure used to achieve these astonishing performances. Important are two types of outputs.

Left: different output neurons represent different actions.

Right: an additional output neuron represents the value of the present state; we can loosely define the value as the probability to win, or the 'average reward' that you can get starting from this state.

The input is a representation of the present state of the game.

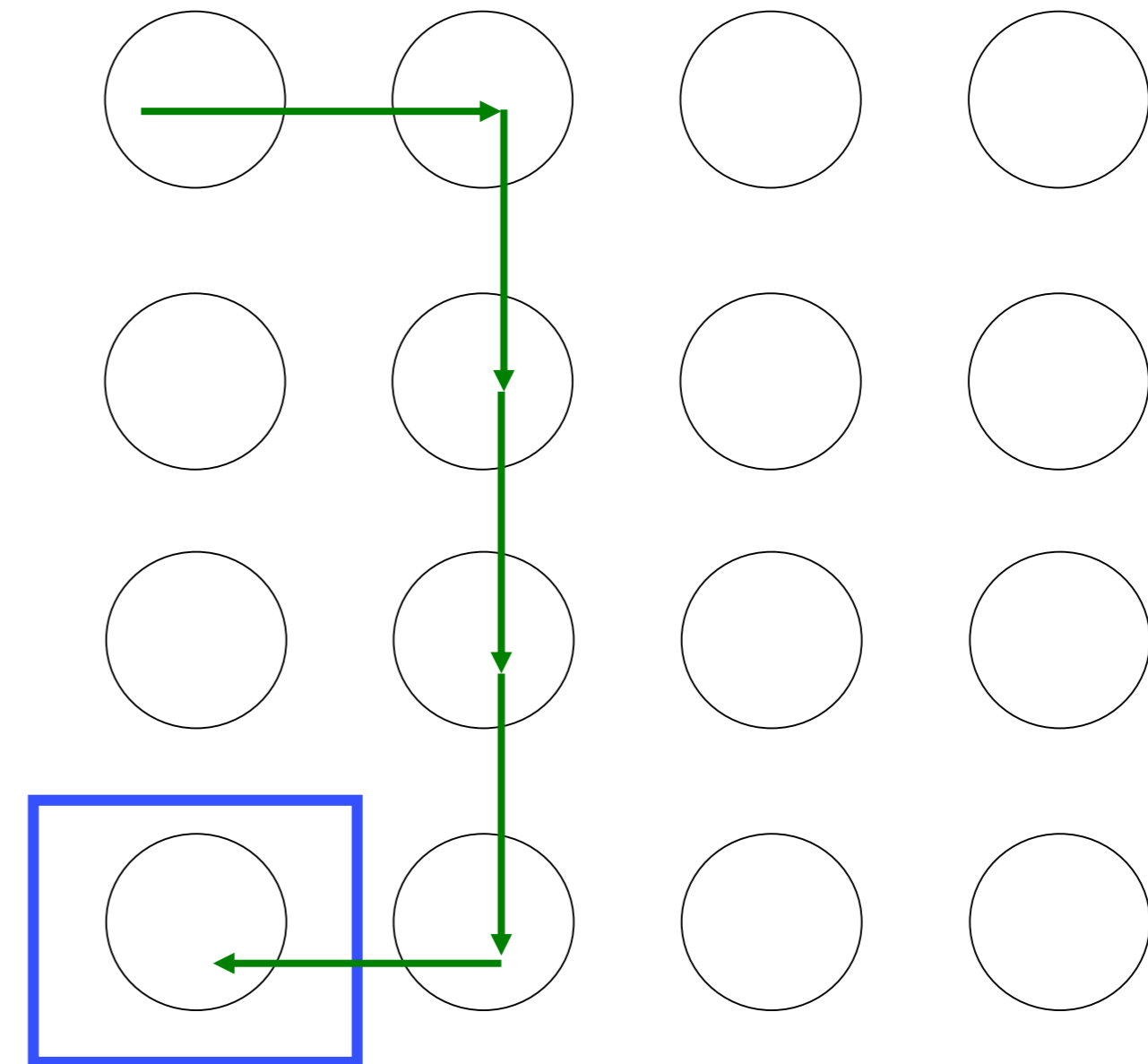
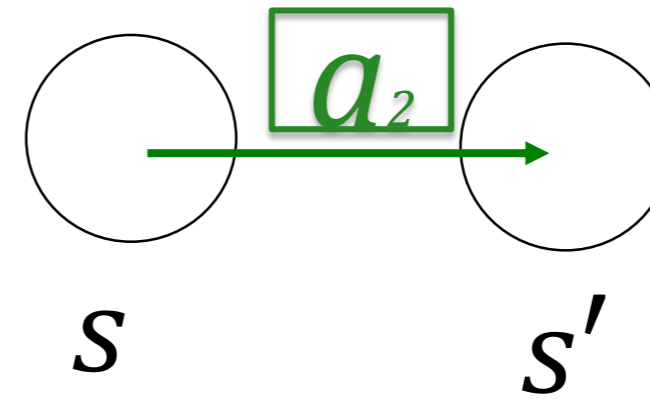
Details will become clear toward the end of the semester; at the moment the aim is just to give you a flavor of the high-level concepts.

# Recall: Elements of Reinforcement Learning:

- discrete states:
  - old state  $s$
  - new state  $s'$
- current state:  $s_t$
- discrete actions:  $a_1, a_2 \dots a_A$
- current action:  $a_t$
- current reward:  $r_t$
- Mean rewards for transitions:

$$R_{s \rightarrow s'}^a$$

often most transitions have zero reward



Previous slide.

The elementary step is:

The agent starts in state  $s$ .

It takes action  $a$

It arrives in a new state  $s'$

Potentially receiving reward  $r$  (during the transition or upon arrival at  $s'$ ).

Since rewards are stochastic we have to distinguish the mean reward at the transition (capital  $R$  with indices identifying the transition) from the actual reward (lower-case  $r$  with index  $t$ ) that is received at time  $t$  on a transition.

Note that in many practical situations most transitions or states have zero rewards, except a single 'goal' state at the end.



# Recall: Elements of Reinforcement Learning

- Environment with discrete states
- Transitions (potentially stochastic) are driven by actions

- **Aim: choose good actions to optimize reward**

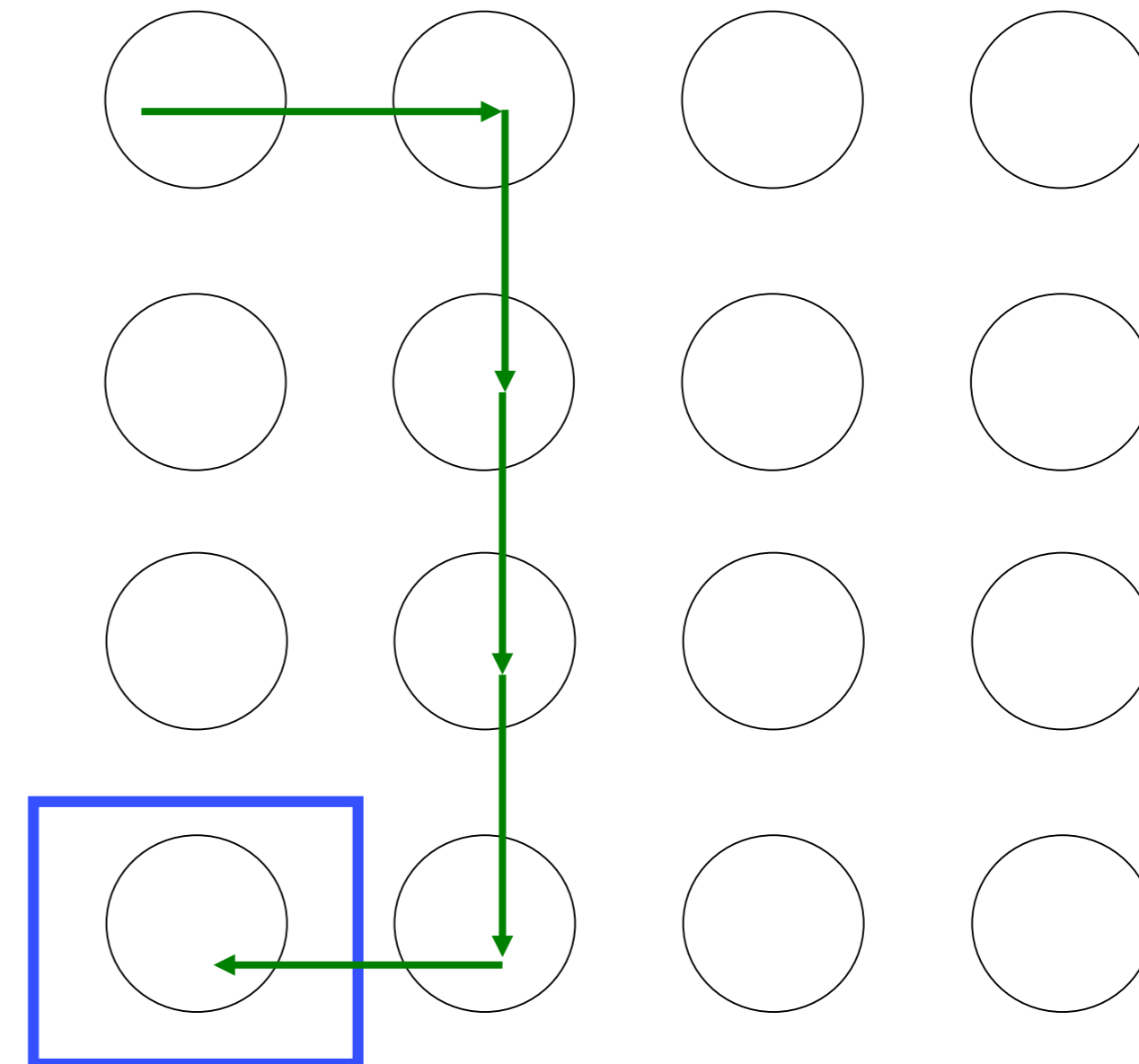
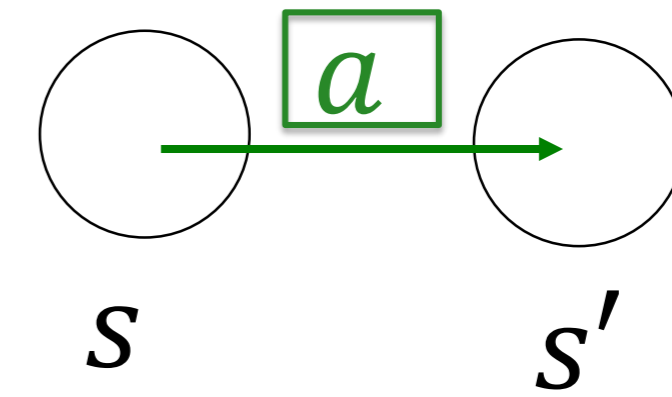
→ **Markov Decision Model**

- Mean reward for transition:

$$R_{s \rightarrow s'}^a = E(r | s, a, s')$$

- **current actual reward:  $r_t$**

often most transitions have zero reward



Previous slide.

Conclusion: In all practical situations, there is an enormous number of states.

In many situations we can think of the actions as discrete.

For the moment we also think of the states as discrete (but next week we will go to continuous state space).

Transitions between actions are influenced by action choices.

Transitions can be stochastic.

Actions should be chosen so as to maximize the reward. This setting is also known as Markov Decision Problem (MDP).

# Recall: Q-value for one-step horizon games/bandit problem

Q-value  $Q(s,a)$

Expected reward for action  $a$  starting from  $s$

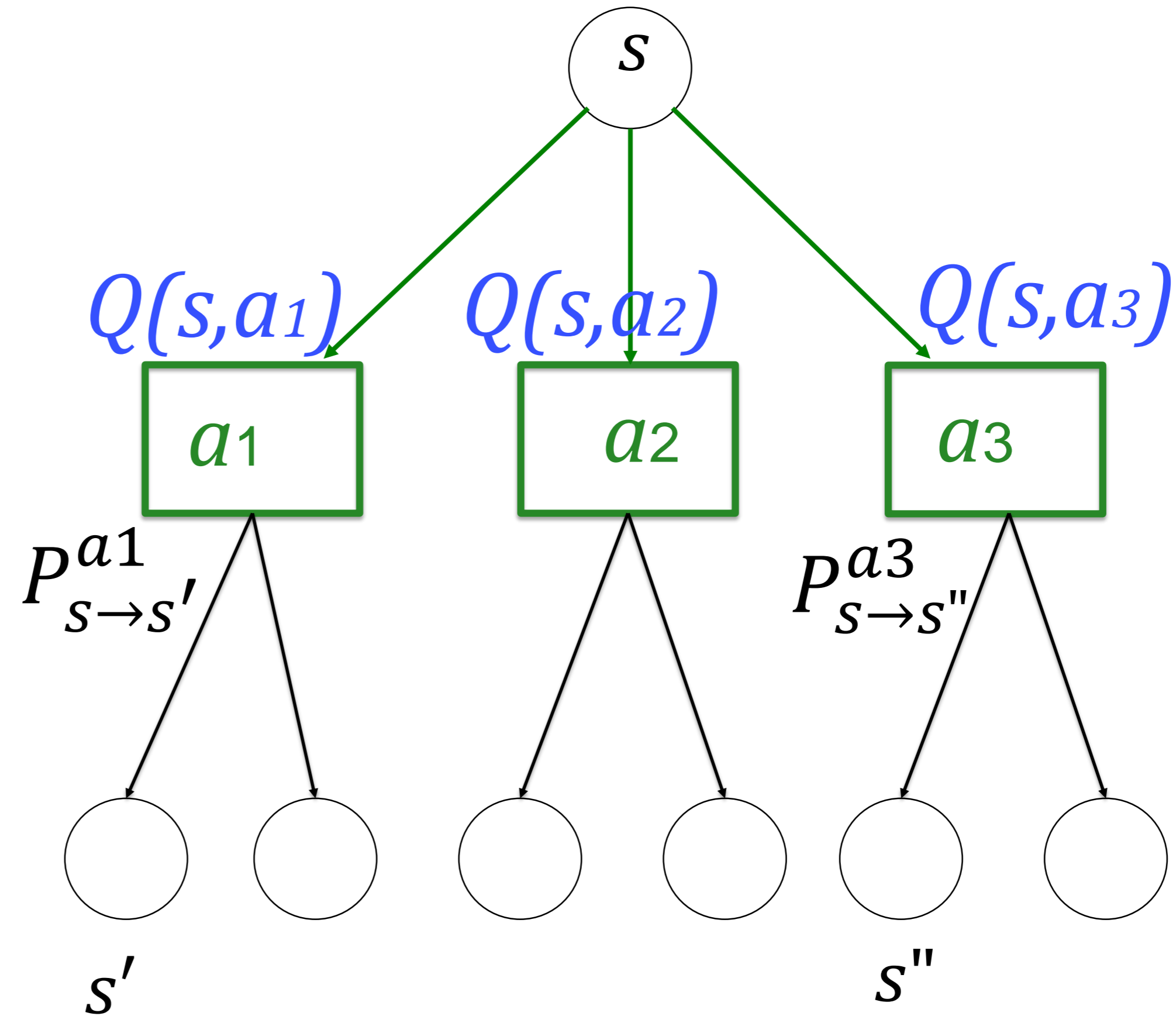
$$Q(s,a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a$$

Reminder:

$$R_{s \rightarrow s'}^a = E(r | s', a, s)$$

Similarly:

$$Q(s,a) = E(r | s, a)$$



Now we know the Q-values: which action should you choose?

Previous slide.

$P_{s \rightarrow s'}^{a1}$  is the probability that you end up in a specific state  $s'$  if you take action  $a1$  in state  $s$ .

We refer to this sometimes as the 'branching ratio' below the 'actions'.

$Q(s,a)$  is attached to the branches linking the state  $s$  with the actions.

actions are indicated by green boxes; states are indicated by black circles.

The mean reward  $R_{s \rightarrow s'}^a$  is defined as the expected reward given that you start in state  $s$  with action  $a$  and end up in state  $s'$  (see Blackboard 1).

Given the branching ratio and the mean rewards, it is easy to calculate the Q-values (Blackboard 1).

# Recall: One-step horizon games (bandit problem)

Q-value = expected reward for state-action pair

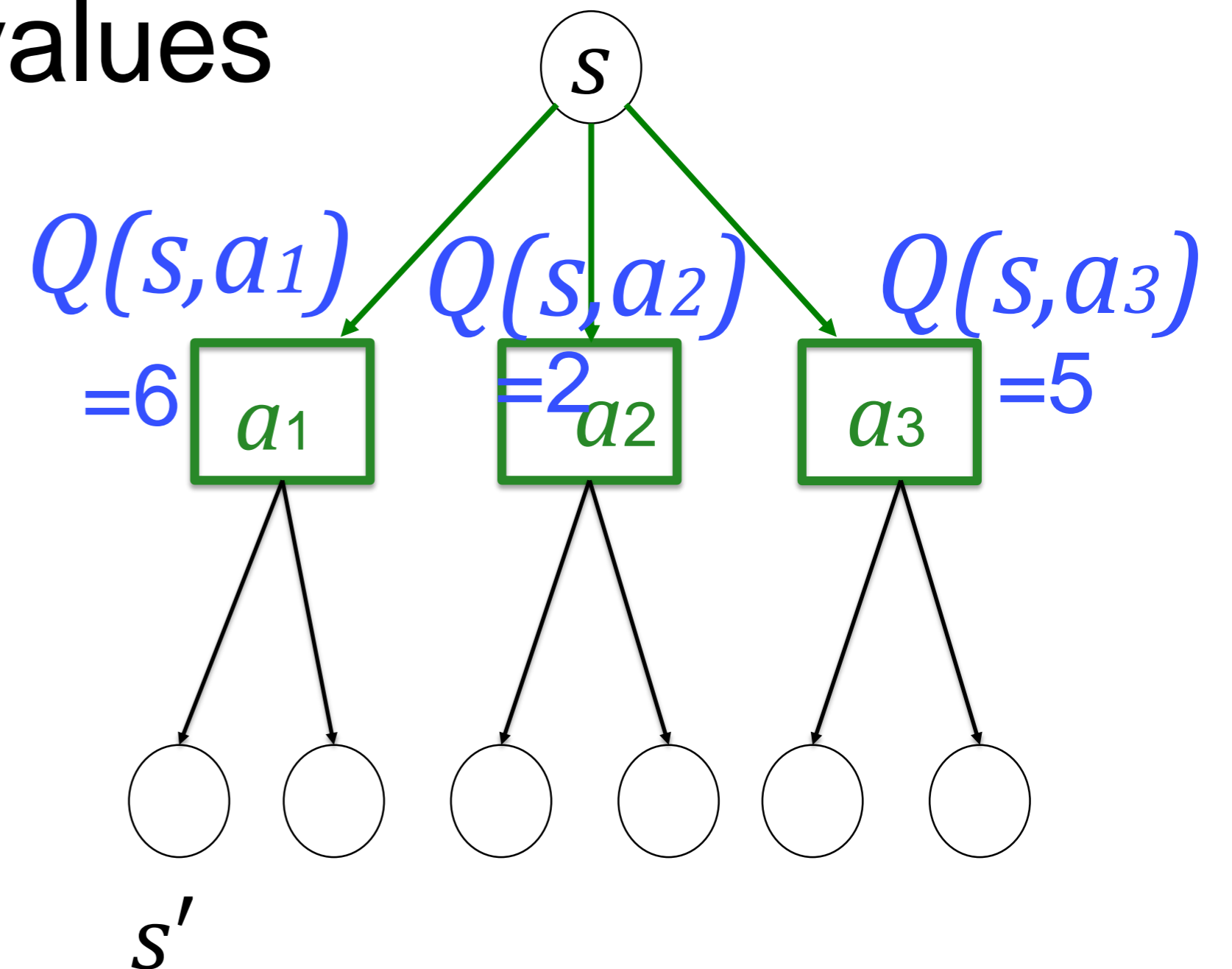
If Q-value is known, choice of action is simple  
→ take action with highest Q-value

BUT: we normally do not know the Q-values  
→ estimate by online update

$$\Delta Q(s, a) = \eta [r_t - Q(s, a)]$$

Learning rate: often  $\eta$ , often  $\alpha$

$$\Delta Q(s, a) = \alpha [r_t - Q(s, a)]$$



Previous slide.

The only remaining problem is that we do not know the Q-values, because the casino gives you neither the branching ratio nor the reward scheme.

Hence the only way to find out is by trial and error (that is, by playing many times – the casino will love this!).

# Recall: Update rule in Expectation (Theorem)

After taking action  $a$  in state  $s$ , we update with

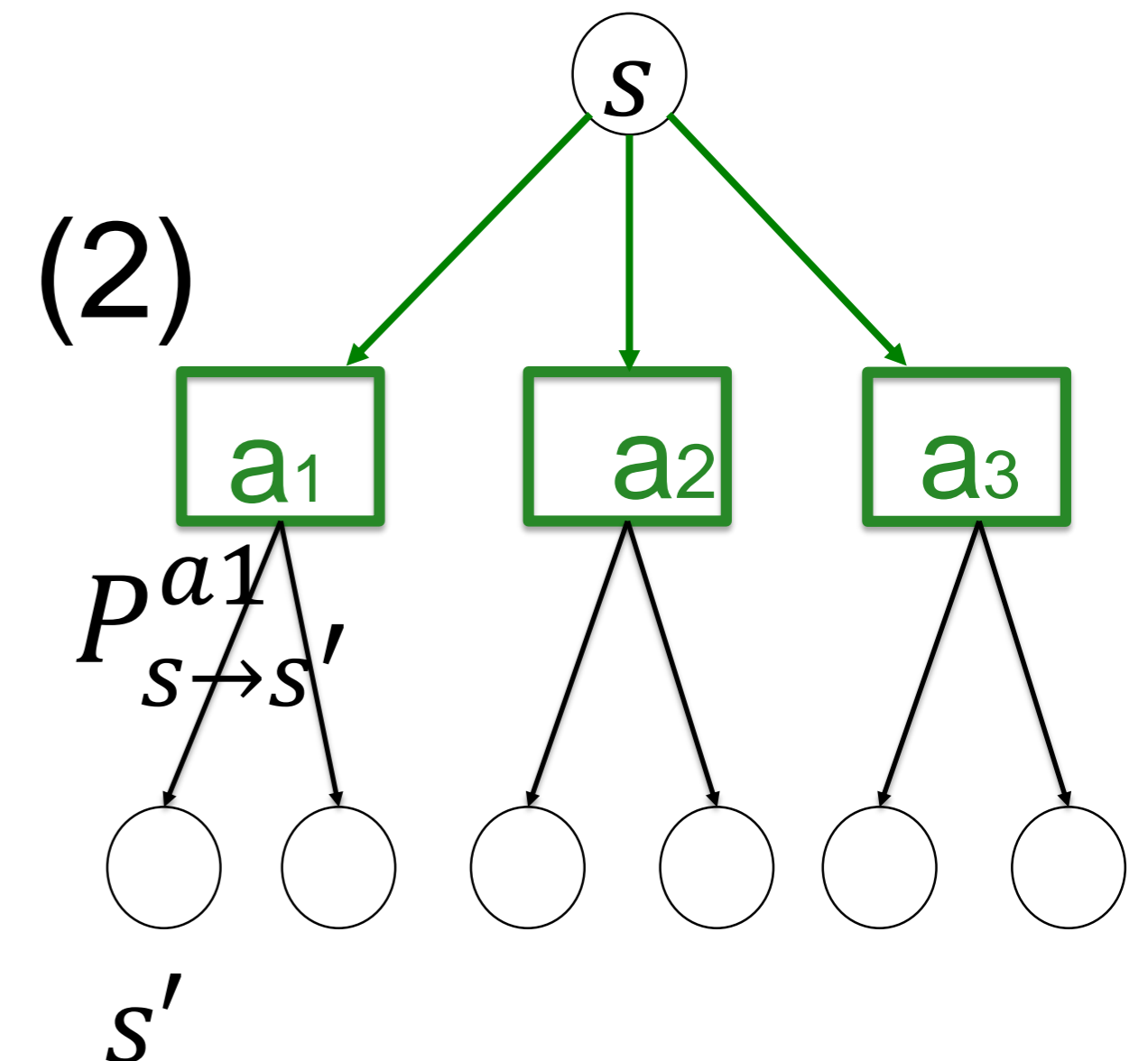
$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)] \quad (1)$$

(i) If the expectation of the update rule (1) **given (s,a)** vanishes, then  $\hat{Q}(s, a)$  has an expectation

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

claim: 'I can even drop the expectation sign in (2)

(ii) If the learning rate  $\eta$  decreases, fluctuations around the **empirical mean**  $\langle \hat{Q}(s, a) \rangle$  decrease and the **empirical mean approaches**  $Q(s, a)$



Previous slide.

When evaluating the **expectation value given  $(s,a)$** , the learning rate drops out since we set the left-hand-side to zero. The exact value of  $\eta$  is not relevant, as discussed in the theorem. Part (i) of the theorem states that the expectation value of  $\hat{Q}(s, a)$  is the correct Q-value. For a quick proof of part (i) see the video. On the blackboard a stronger statement was shown. Today we make the statements more precise. Convergence in expectation is equivalent to imagining that you start millions of trials with the same value  $\hat{Q}(s, a)$  without any intermediate update. So in that sense it is like a super-big 'batch' of examples.

**In practice, we do not have expectations but online updates with fluctuations.** It is important is that  $\eta$  is small at the end of learning so as to limit the amount of fluctuations. Part (ii) states that **online mean** for small learning rate also goes to the correct Q-value. Indeed, since the equations are linear (for the bandit problem = 1-step horizon), the calculation of part (i) apply analogously to the long-term empirical temporal average (denoted by angular brackets)

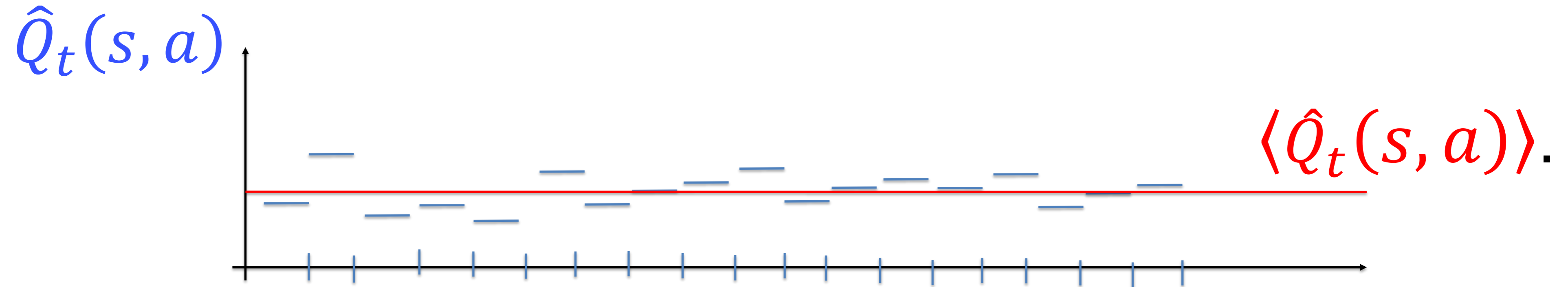
$$\langle \Delta \hat{Q}(s, a) \rangle = \eta \langle [r_t - \hat{Q}(s, a)] \rangle$$

This equivalence based on linearity is not true for the multi-step horizon that we discuss later in this lecture.



Part (ii) of Theorem:

We work with the **online update**  $\Delta \hat{Q}(s, a)$  . With finite learning rate, the value of  $\hat{Q}_t(s, a)$  fluctuates around a mean  $\langle \hat{Q}_t(s, a) \rangle$ .



**Claim:** Under the hypothesis  $\langle \Delta \hat{Q}(s, a) = 0 \rangle$  , the mean  $\langle \hat{Q}_t(s, a) \rangle$  is equal to the ‘correct’ Q-value.

Your notes. (Proof in the Blackboard notes)

# One-step horizon: summary

Q-value = expected reward for state-action pair

If Q-value is known, choice of action is simple

→ take action with highest Q-value

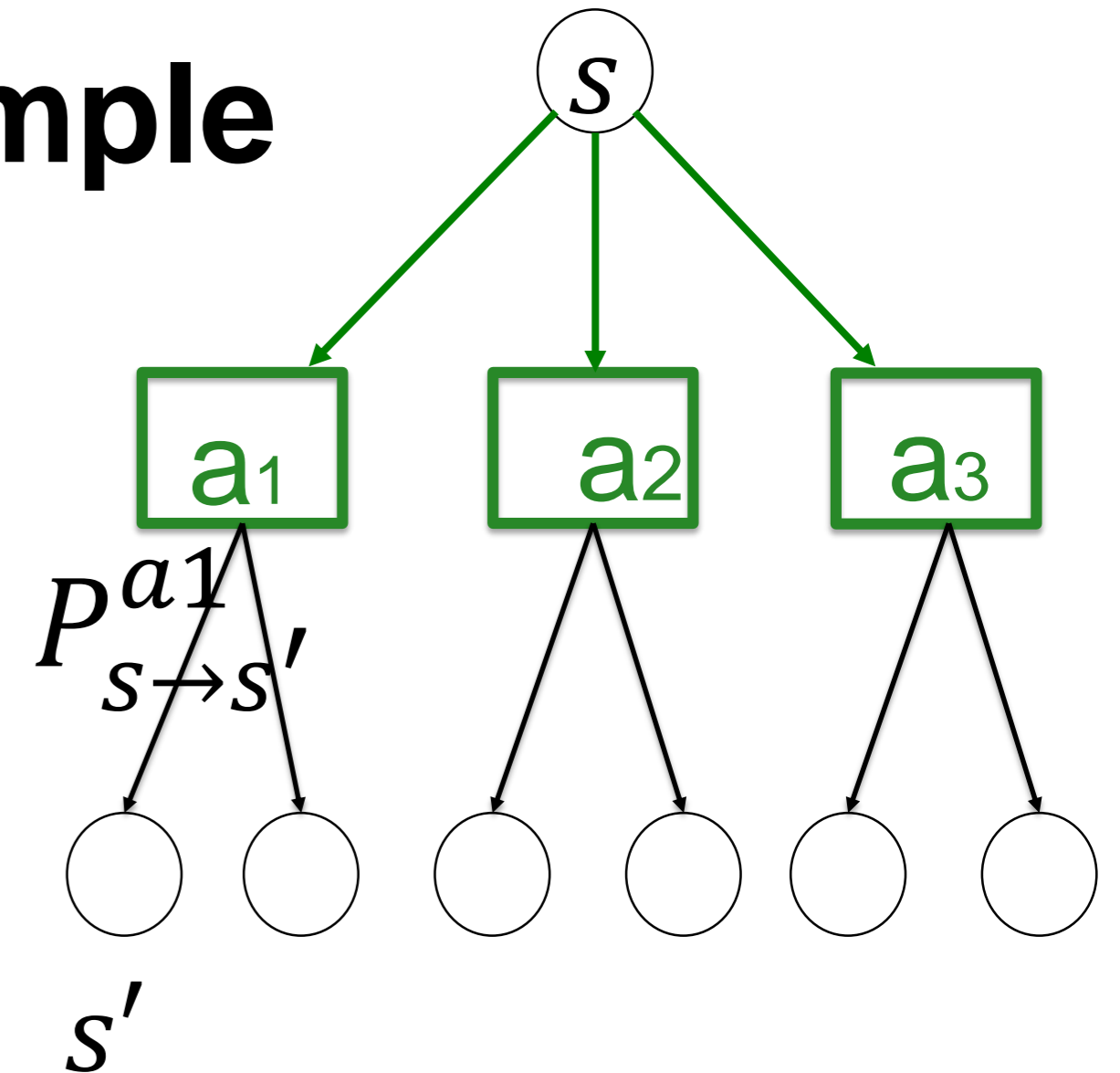
If Q-value not known:

→ estimate  $\hat{Q}$  by trial and error

→ update with rule

$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)] \quad (1)$$

→ Let learning rate  $\eta$  decrease over time



Iterative algorithm (1) fluctuates, but ‘makes sense’

Previous slide.

Let us distinguish the ESTIMATE  $\hat{Q}(s, a)$  from the real Q-value  $Q(s, a)$

The update rule can be interpreted as follows:

if the actual reward is larger than (my estimate of) the expected reward, then I should increase (a little bit) my expectations.

The learning rate  $\eta$  :

In exercise 1, we found a rather specific scheme for how to reduce the learning rate over time. But many other schemes also work in practice. For example you keep  $\eta$  constant for a block of time, and then you decrease it for the next block.

**Note: in later lectures I will often use the symbol  $\alpha$  instead of  $\eta$**

Both symbols indicate what is called the 'learning rate' in Deep Learning.

# Artificial Neural Networks: RL1

## Reinforcement Learning and SARSA

Wulfram Gerstner

EPFL, Lausanne, Switzerland

Part 3b: Expectation, Batch, and ONLINE rules

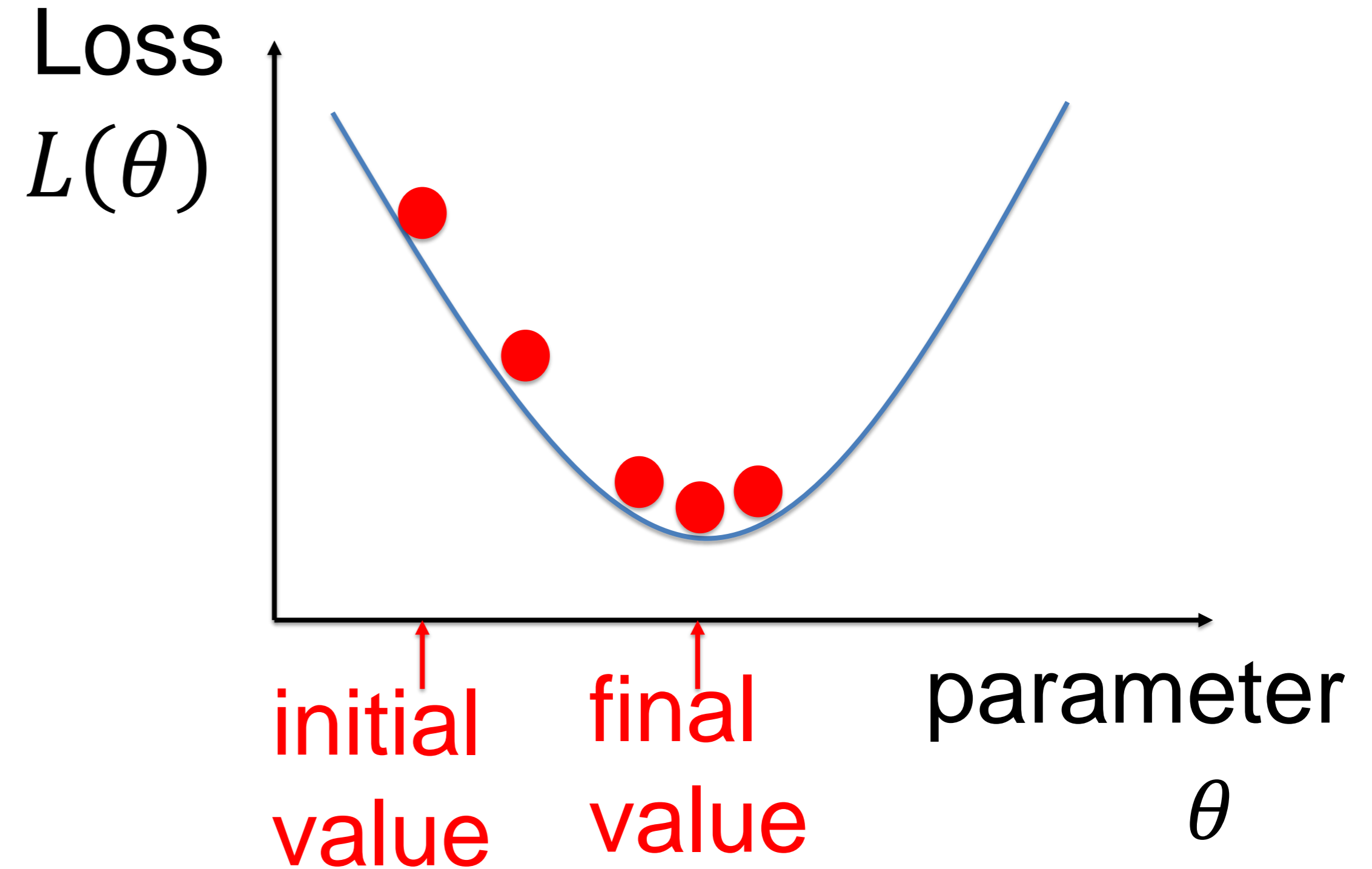
- Examples of Reward-based Learning
- Elements of Reinforcement Learning
- One-step horizon (bandit problems)
- **Expectation, batch, and online rules**

Previous slide.

Last week we did a first calculation with expectations and I argued that this is 'batch-like'. Since the type of calculation is important, but since this comparison caused many questions after class, I add a detour.

All the material in this part is in principle standard material in Machine Learning classes, even though I put the accent on aspects that are important for Reinforcement Learning.

# Detour Machine Learning ReCap: Online, Batch, Expectation



Your notes. (Review of gradient Descent)

The set of parameters (also called parameter vector) is generically denoted by  $\theta$ .

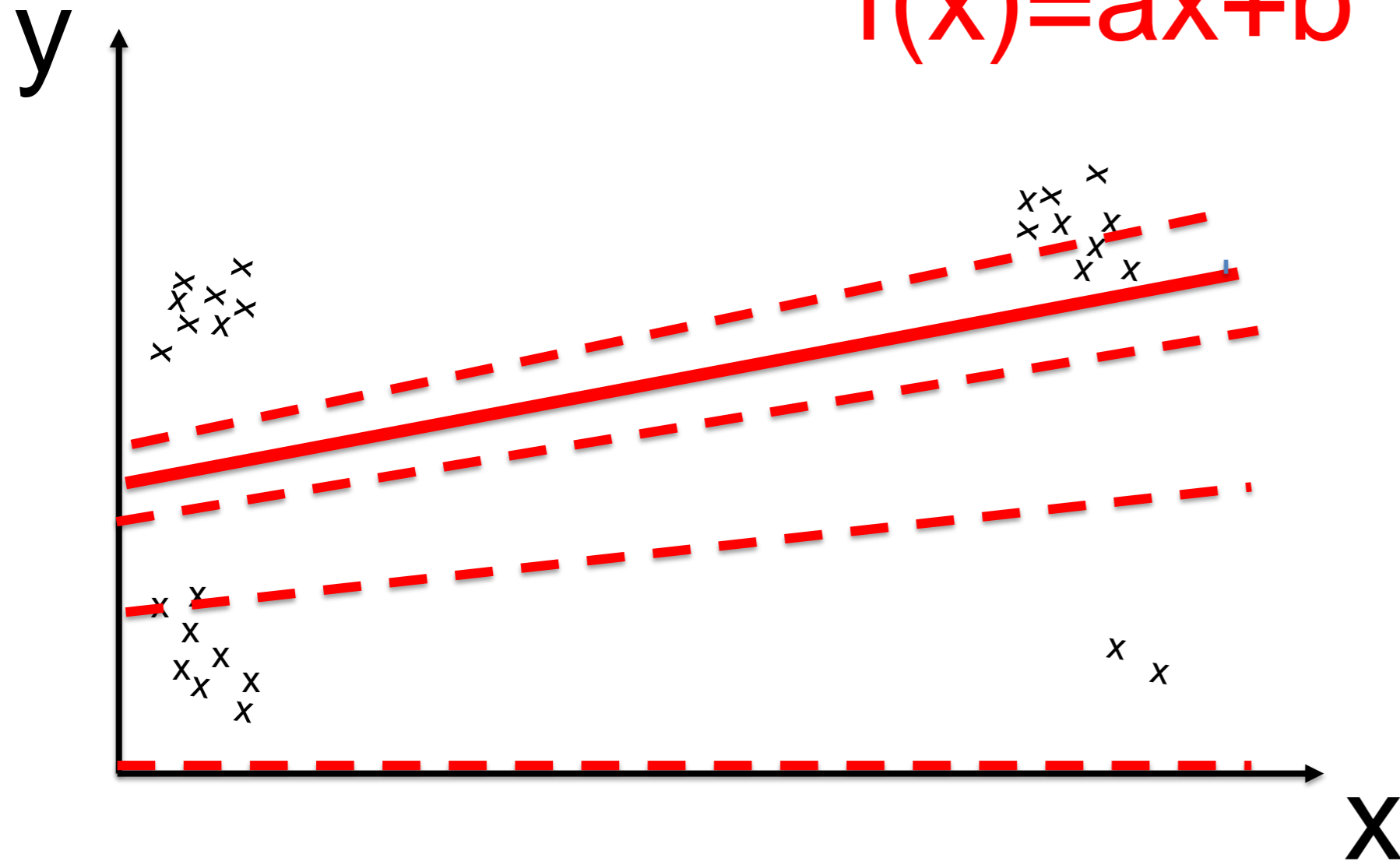
The loss function (error function) is denoted by capital  $L$ .

.

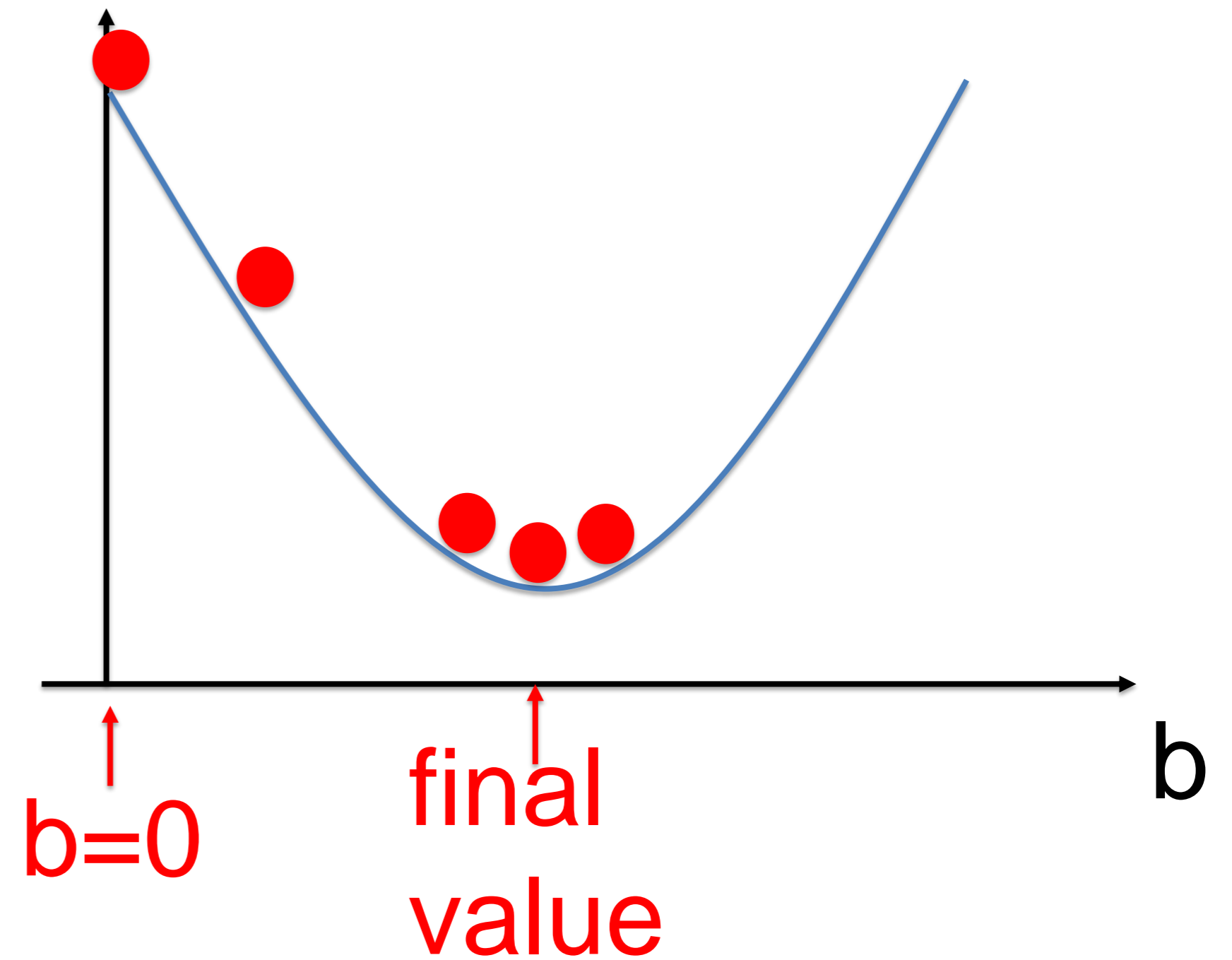


# Detour Machine Learning ReCap: Online, Batch, Expectation

$$f(x) = ax + b$$



Loss



Gradient descent in 'batch mode'

Your notes. (Review of gradient Descent)

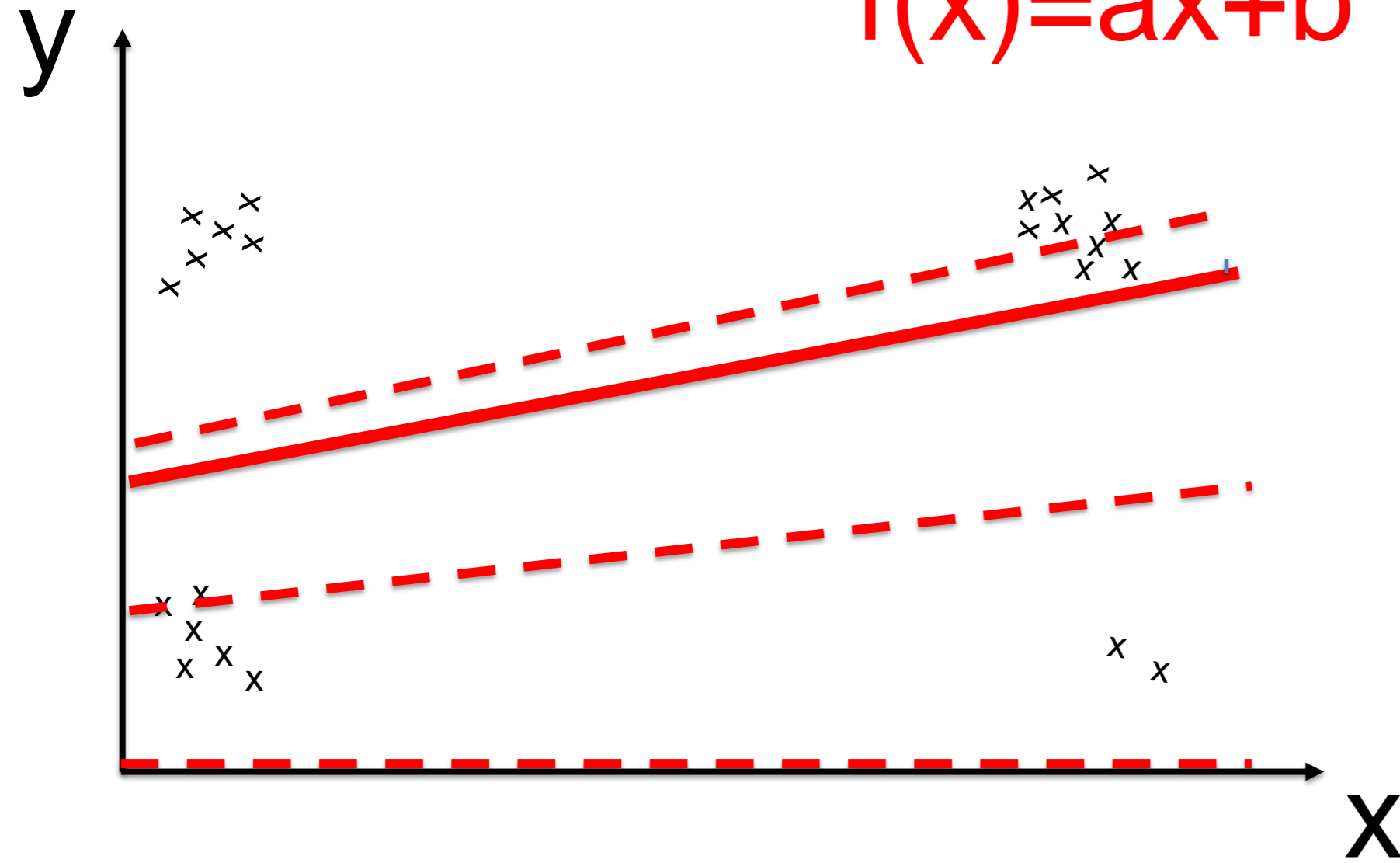
The specific parameters of the linear function are  $a$  and  $b$ . For the drawing of the loss function, only one of the two parameters is plotted.

The aim is to fit a set of data points by the linear function.  
Dashed red lines show intermediate update steps.

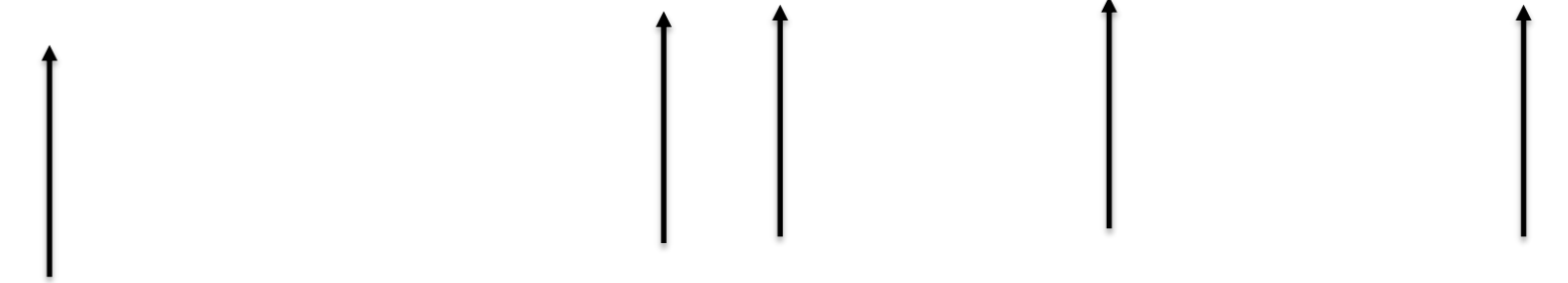
**Note: For the learning rate I will often use the symbol  $\alpha$  instead of  $\eta$**

# Detour Machine Learning ReCap: Online, Batch, Expectation

$$f(x) = ax + b$$



$$f(x|\theta) = f(x|a, b) = ax + b$$



Parameters  $\theta = (a, b)$

$$\Delta\theta = -\alpha \frac{d}{d\theta} L(\theta)$$

$$\Delta b = -\alpha \frac{\partial}{\partial b} L(a, b)$$

algo (iterative update)  
*iterate to convergence criteria*

$$b^{old} \leftarrow b$$

$$b = b^{old} + \Delta b$$

*analogously for a*

Your notes. (Review of gradient Descent)

The set of parameters (also called parameter vector) is generically denoted by  $\theta$ .  
And then the parameters are specified to be  $a$  and  $b$ .

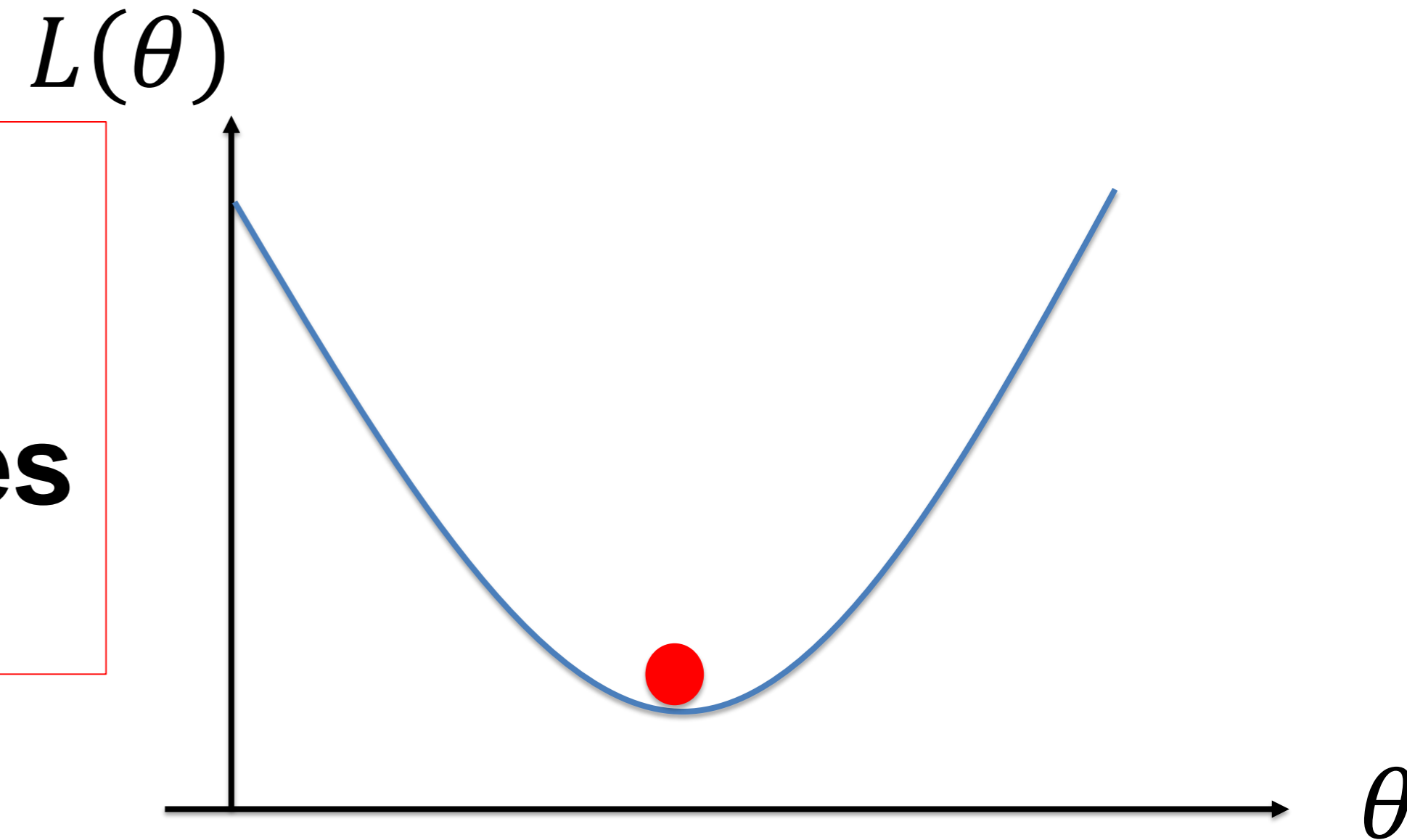
Changes of parameters are calculated by gradient descent on the Loss function.

# Detour Machine Learning ReCap: Online, Batch, Expectation

**CONCLUSION 1, from rule (1):**

**if  $\Delta\theta=0$  then**

- **parameter  $\theta$  no longer changes**
- **(local) minimum at  $\theta^{old}$**



Gradient is always evaluated at  $\theta^{old}$

$$\Delta\theta = -\alpha \frac{d}{d\theta} L(\theta) \Big|_{\theta^{old}} \quad (1)$$

algo (iterative update)  
*iterate to convergence criteria*

$$\theta^{old} \leftarrow \theta$$

$$\theta = \theta^{old} + \Delta\theta$$

Your notes. (Review of gradient Descent)

The update rule tells us immediately that the update vanishes at a parameter value that has zero gradient. Only minima can be generically approached by gradient descent (not the maxima).

# Detour Machine Learning ReCap: Online, Batch, Expectation

Loss function

$$L(\theta) = \mathbf{E}[l(f(x|\theta), y)]$$

$$L(\theta) = \frac{1}{N} \sum_k^N [l(f(x_k|\theta), y_k)]$$

loss per data point

Gradient descent (batch)

$$\Delta\theta = -\alpha \frac{d}{d\theta} L(\theta)$$

Example:  $l = \text{L2 loss}$ , linear model

$$f(x_k|\theta) = f(x_k|a, b) = ax_k + b$$

$$L(\theta) = \frac{1}{N} \sum_k^N (ax_k + b - x_k)^2$$

$$\Delta\theta = -\alpha \mathbf{E}\left[\frac{d}{d\theta} l(f(x|\theta), y)\right]$$

$$\Delta\theta = -\alpha \frac{1}{N} \sum_k^N \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]$$

Your notes. (Review of gradient Descent)

The loss function is the expectation across all possible pairs  $(x,y)$  with the appropriate statistical weight. The loss per data point is denote by a small character  $l$ .

Often a large batch of  $N$  data points is taken instead. These  $N$  data points must be representative for the statistical distribution  $p(x,y)$ .

The example shows the linear function that we considered earlier.



# Detour Machine Learning ReCap: Online, Batch, Expectation

Loss function

$$L(\theta) = \mathbf{E}[l(f(\mathbf{x}|\theta), y)]$$

$$L(\theta) = \frac{1}{N} \sum_k^N [l(f(x_k|\theta), y_k)]$$

loss per data point

$$L(\theta) = \int dx dy p(x, y) [l(f(\mathbf{x}|\theta), y)]$$

Gradient descent (batch)

$$\Delta\theta = -\alpha \frac{d}{d\theta} L(\theta)$$

$$\Delta\theta = -\alpha \int dx dy p(x, y) \left[ \frac{d}{d\theta} l(f(x|\theta), y) \right]$$

$$\Delta\theta = -\alpha \mathbf{E} \left[ \frac{d}{d\theta} [l(f(x|\theta), y)] \right]$$

$$\Delta\theta = -\alpha \frac{1}{N} \sum_k^N \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]$$

Your notes. (Review of gradient Descent)

On the left:

The loss function is the expectation across all possible pairs  $(x,y)$  with the appropriate statistical weight. The gradient operation is linear and can be exchanged with the expectation (which is also a linear operation).

On the right:

The same calculation with a large batch of  $N$  data points.

The average of  $N$  in the gradient is analogous to the expectation (red boxes).

# Conclusion: Expectation = Batch size N to infinity

$$\Delta\theta = -\alpha E\left[\frac{d}{d\theta} l(f(x|\theta), y)\right]$$

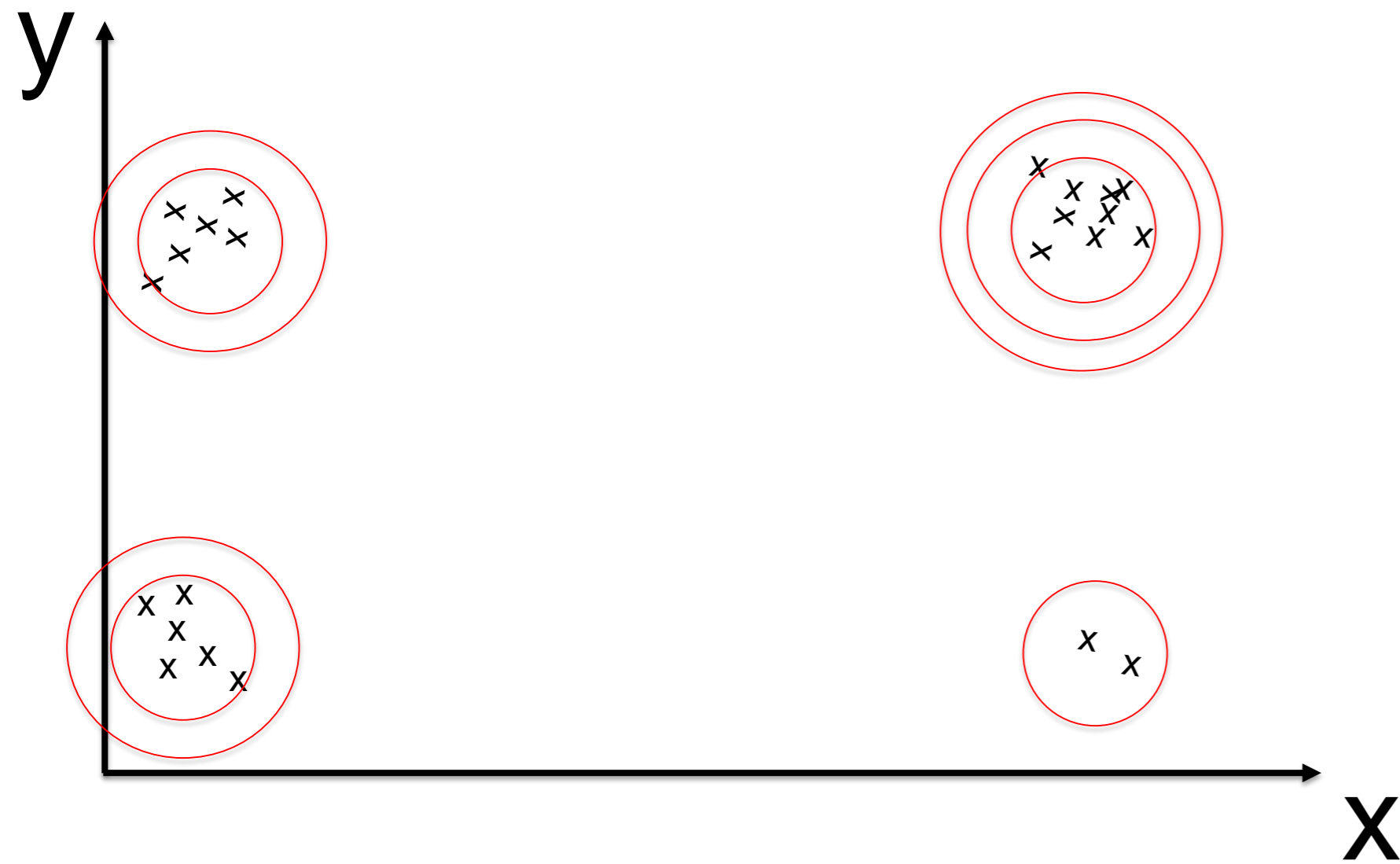
$$\Delta\theta = -\alpha \int dx dy p(x, y) \left[\frac{d}{d\theta} l(f(x|\theta), y)\right]$$

‘statistical  
formulation’  
with ‘expectations’

choose N data points using the  
appropriate statistical weight

$$\Delta\theta = -\alpha \frac{1}{N} \sum_k^N \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]$$

$$E[\dots] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_k^N [\dots]$$



Your notes. (Review of gradient Descent)

We said that the average of  $N$  in the gradient is 'analogous' to the expectation. For the limit  $N$  to infinity batch and expectation are again identical.

The idea of the density  $p(x,y)$  is shown for the same example as before.

# Conclusion: Expectation = Batch size N to infinity

$$\Delta\theta = -\alpha E\left[\frac{d}{d\theta} l(f(x|\theta), y)\right]$$

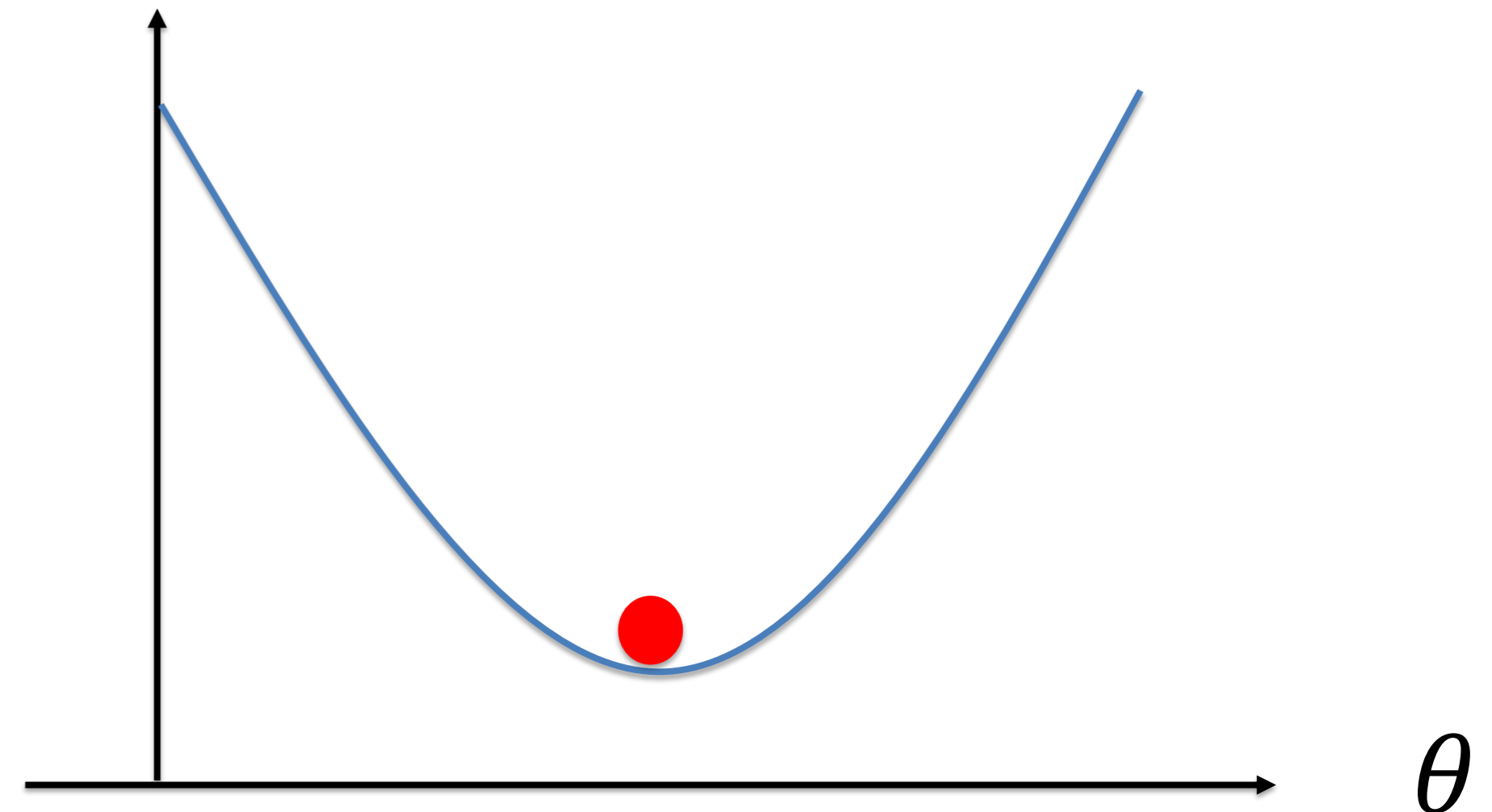
$$E[\dots] = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_k^N [\dots]$$

$$\Delta\theta = -\alpha \frac{1}{N} \sum_k^N \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]$$

**CONCLUSION 1 from rule (1):**

if  $\Delta\theta=0$  (with N to infinity) then

- $\theta$  doesn't change
- (local) minimum at  $\theta^{old}$
- $\theta^{old} = \theta^{optim}$  in 'statistical' sense



Your notes. (Review of gradient Descent)

A repetition of what we have seen before:

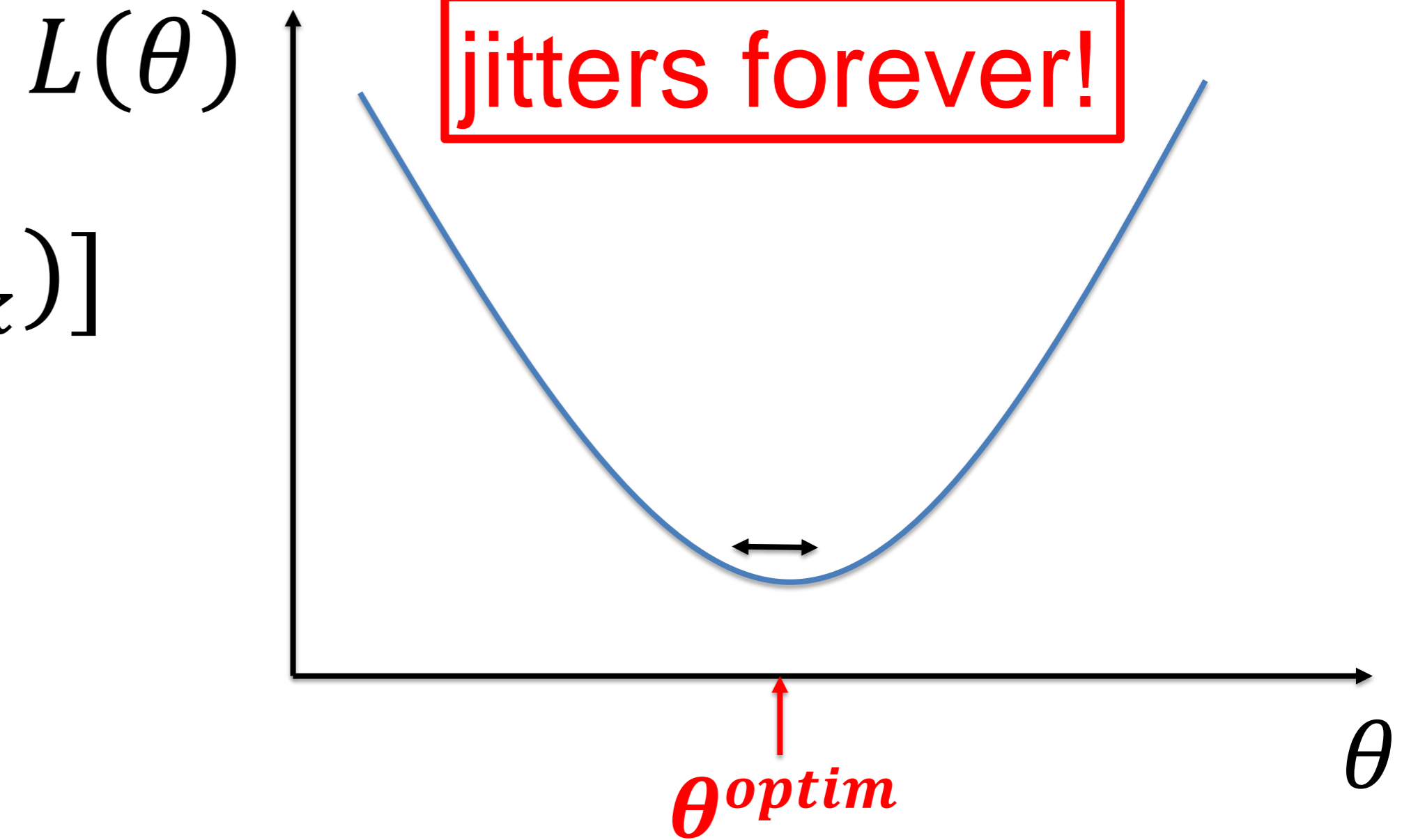
If the update step in the batch rule (N to infinity) vanishes, then we know that we are at a minimum of the loss function.

And this is equivalent to saying:

If the update step of the true loss function with the expectation sign vanishes, then we know that we are at a minimum of the loss function.

# Detour Machine Learning ReCap: Batch versus 'Online'

$$\Delta\theta = -\alpha \frac{1}{N} \sum_k^N \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]$$



**Online:**

$$\Delta\theta = -\alpha \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]_{\theta=\theta^{old}}$$

Update after each data point

algo (iterative update)  
*iterate to convergence criteria*

$$\theta^{old} \leftarrow \theta$$

$$\theta = \theta^{old} + \Delta\theta$$

Your notes. (Review of gradient Descent)

An online rule means: drop the statistical averaging.

Here it means: drop the sum over data points.

As a result the parameter vector  $\theta$   
can change after each data point!

And this is true even if we are already at the exact minimum of the true loss. The next data point might for example be an outlier and the parameter vector changes again. Therefore the gradient descent solution always jitters.

The size of the jitter depends on the learning rate (here called alpha).



# Detour Machine Learning ReCap: Batch, 'Online', Expectation

**Online:**

$$\Delta\theta = -\alpha \frac{d}{d\theta} [l(f(x_k|\theta), y_k)]_{\theta=\theta^{old}}$$

Update after each data point

**Conclusion:**

- Online update has jitter

**BUT**

- Expected update has no jitter

**Expected Online Update ( $\theta = \theta^{old}$  frozen):**

$$E[\Delta\theta] = -\alpha E\left[\frac{d}{d\theta} [l(f(x_k|\theta), y_k)]_{\theta=\theta^{old}}\right]$$

**Conclusion:**

Expected update of the online rule is identical to batch update with infinite data

Your notes. (Review of gradient Descent)

We can ask:

What would be the EXPECTATION of the update step.

Suppose we momentarily have the parameter  $\theta = \theta^{old}$ .

Then we ask what is the EXPECTED change at this location.

Comparison with the batch rule shows that the expected update of the online rule is identical to batch update with infinite data evaluated at  $\theta = \theta^{old}$ .

THIS IMPLIES:

**If by chance  $\theta = \theta^{old}$  is the exact minimum, the expected update is zero; but the ACTUAL update can be nonzero!**

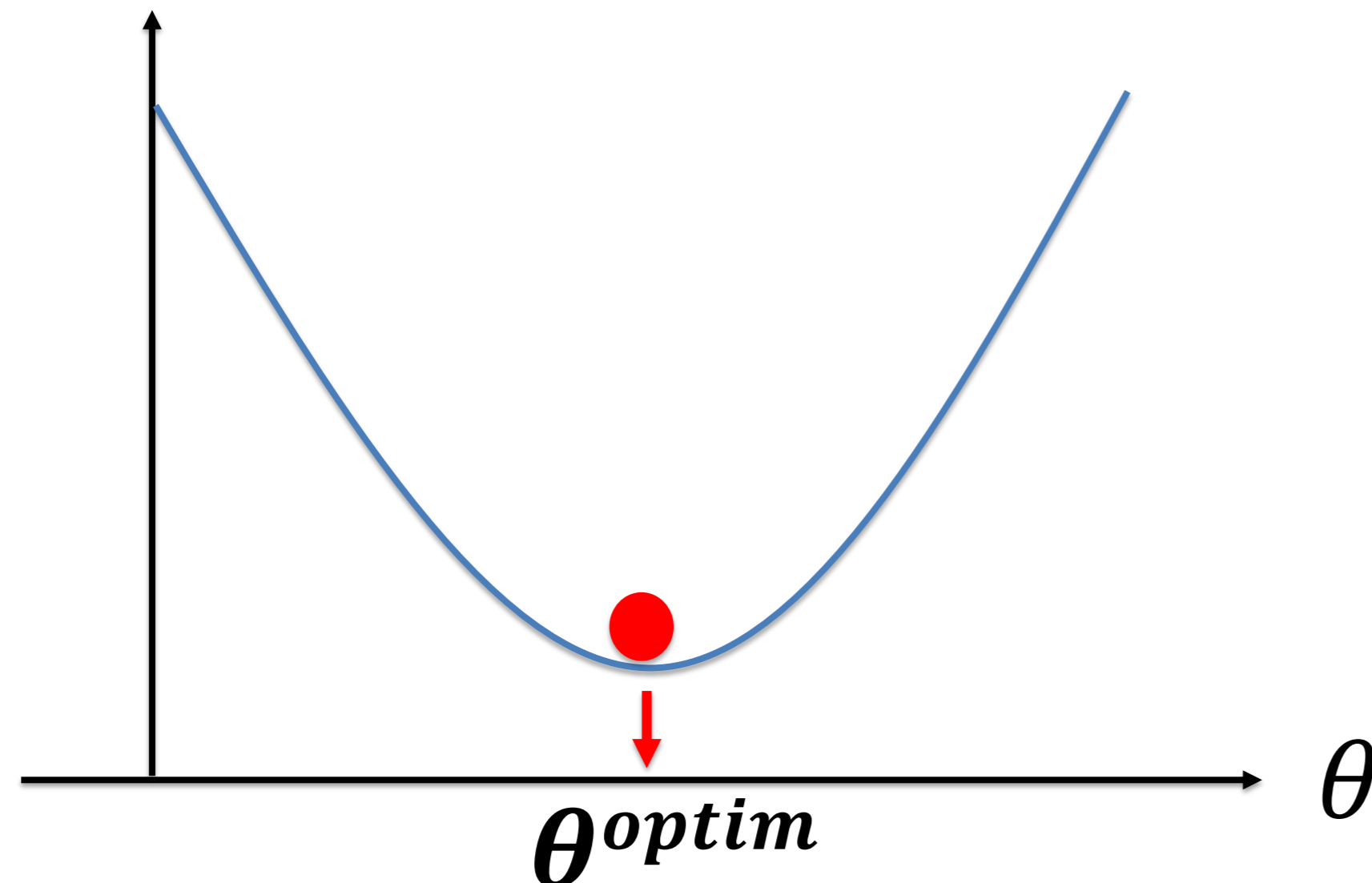
This is also summarized in the next slide and the quiz.

## Batch rule with N to infinity

**CONCLUSION 1** from rule (1):

if  $\Delta\theta=0$  (with N to infinity) then

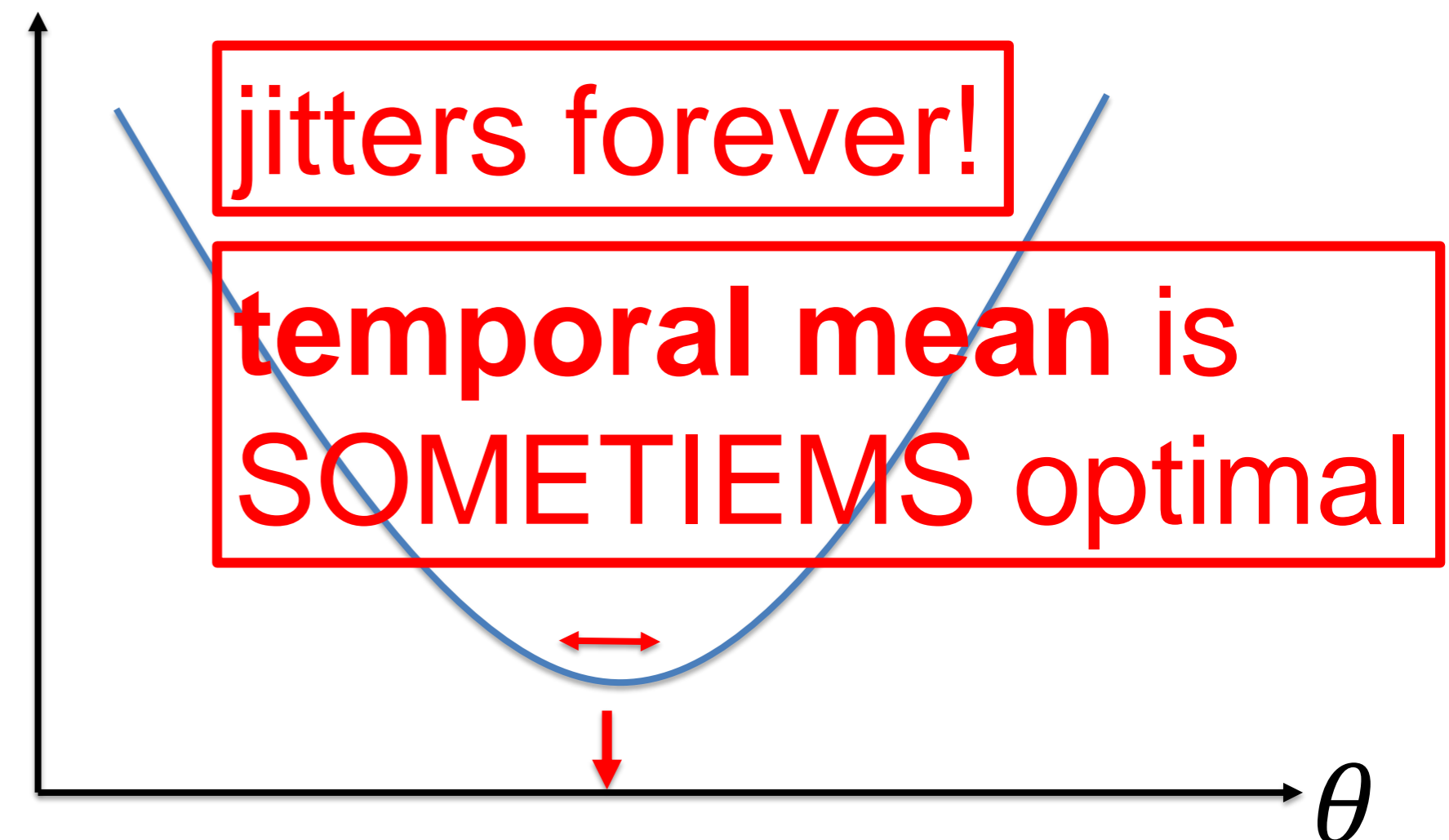
- $\theta$  doesn't change
- (local) minimum at  $\theta^{old}$
- $\theta^{old} = \theta^{optim}$



## Online Rule

**CONCLUSION 2:**

- $\theta$  jitters forever. **BUT:**
- if by chance  $\theta^{old}$  such that  $E(\Delta\theta)=0$  then (local) minimum at  $\theta^{old} = \theta^{optim}$



Previous slide/next slide. Summary slides:

If the expected update is zero [i.e.,  $E(\Delta\theta)=0$  ] for a given set of parameter  $\theta = \theta^{old}$ , then  $\theta$  is a locally optimal parameter, **even for the online rule:**

$$\theta = \theta^{old} = \theta^{optim}$$

There is no statement how we would find this parameter  $\theta = \theta^{optim}$

A completely different statement concerns the **mean of the jittering parameter  $\theta$**  .

If the **update steps are symmetric**, then the mean of the parameter  $\theta$  is the optimal one:

$$\langle \theta \rangle = \theta^{optim}$$

However, if the update steps are asymmetric, then the mean  $\langle \theta \rangle$  of the parameter  $\theta$  is shifted compared to the optimal one (next slide).

For gradient descent on a loss function we recognize the asymmetry in loss curve.

However (even in cases where we do not have a loss function) what really counts is whether the update steps are symmetric or not:

Suppose the current parameter is  $\theta = \theta^{optim} + \epsilon$  where  $\epsilon$  is small, i.e. close to the optimum

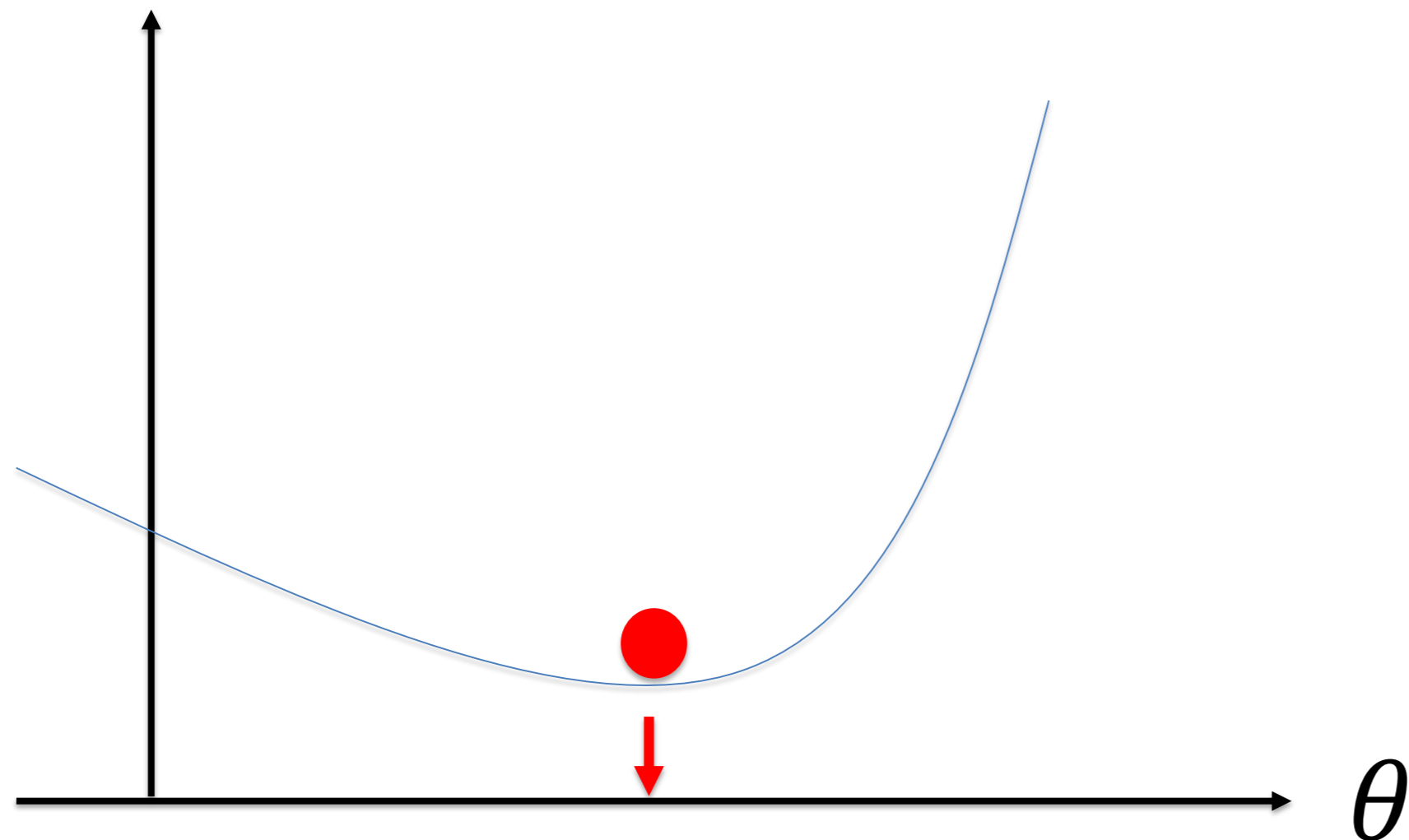
Symmetry is guaranteed if update steps are linear  $\Delta\theta = \alpha\epsilon$  with small constant  $\alpha$ .

## Batch rule with N to infinity

**CONCLUSION 1** from rule (1):

if  $\Delta\theta=0$  (with N to infinity) then

- $\theta$  doesn't change
- (local) minimum at  $\theta^{old}$
- $\theta^{old} = \theta^{optim}$



## Online Rule

**CONCLUSION 2:**

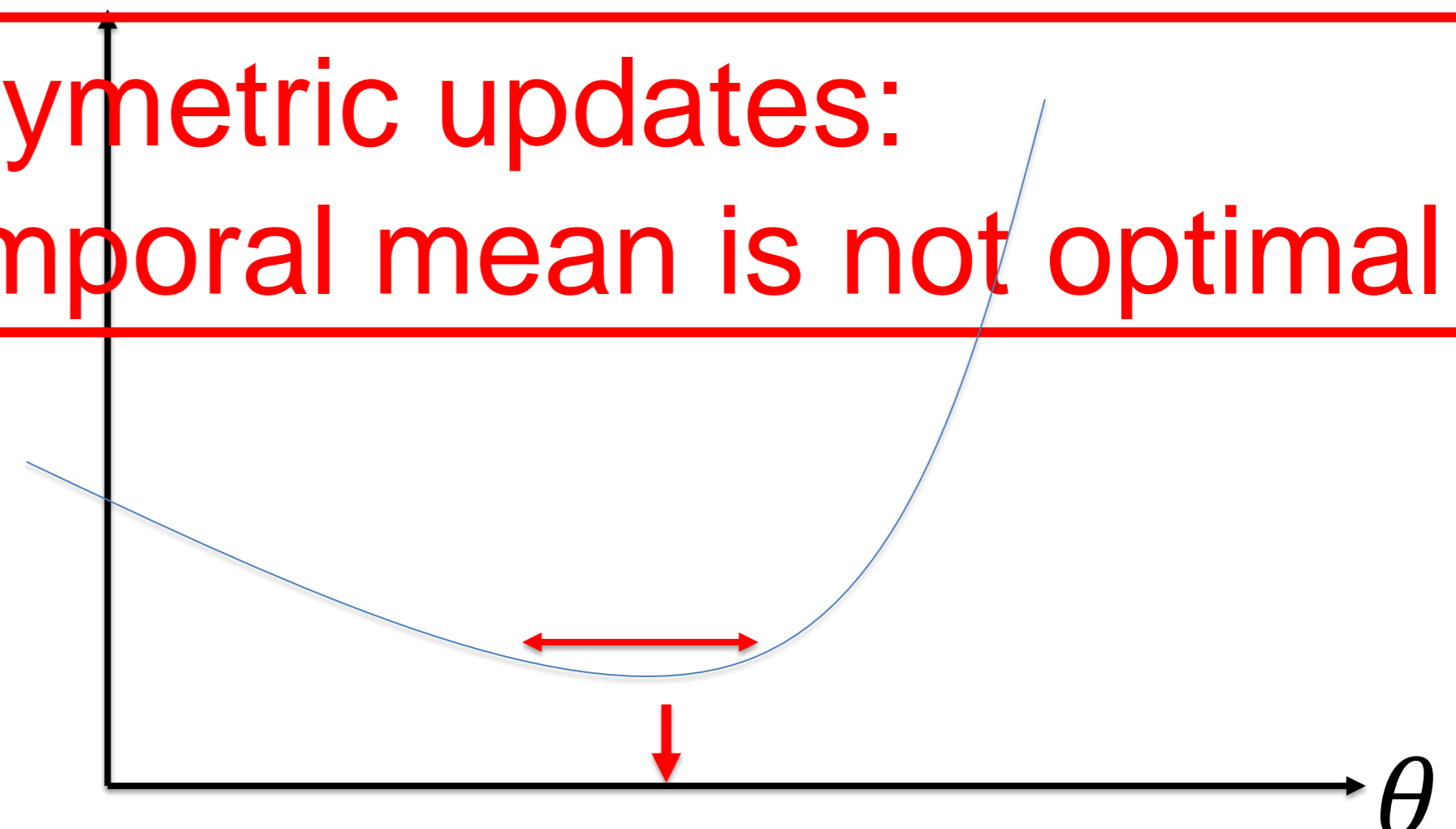
-  $\theta$  jitters forever. **BUT:**

- if by chance  $\theta^{old}$  such that

$E(\Delta\theta)=0$  then

(local) minimum at  $\theta^{old} = \theta^{optim}$

Asymmetric updates:  
temporal mean is not optimal



## Quiz: Expectation, Batch, Online (Recap of ML)

- [ ] With a **batch rule** and small learning rate, I sometimes reach a local minimum without remaining parameter jitter.
- [ ] With a **batch rule at a local minimum** I never have any remaining parameter jitter
- [ ] With an **online rule at a local minimum** I never have any remaining parameter jitter
- [ ] With an **online rule** at a local minimum the **expectation of the online update step vanishes.**
- [ ] The **expectation of the online update step** is equivalent to a very large **batch** (N to infinity)
- [ ] With an online rule jittering round the minimum, the temporal mean is guaranteed to be at the location of the minimum

# Teaching monitoring – monitoring of understanding

up to here, at least 60% of material was new to me.

I have the feeling that I have been able to follow  
(at least) 80% of the lecture up to here.

**End of Detour:**

**Apply to Q-values in the Bandit problem.**

ML: parameters are called  $\theta$

Function fitting: parameters are a and b

Bandit problem: parameters are  $Q(s,a)$



# Recall: Update rule in Expectation (weak version)

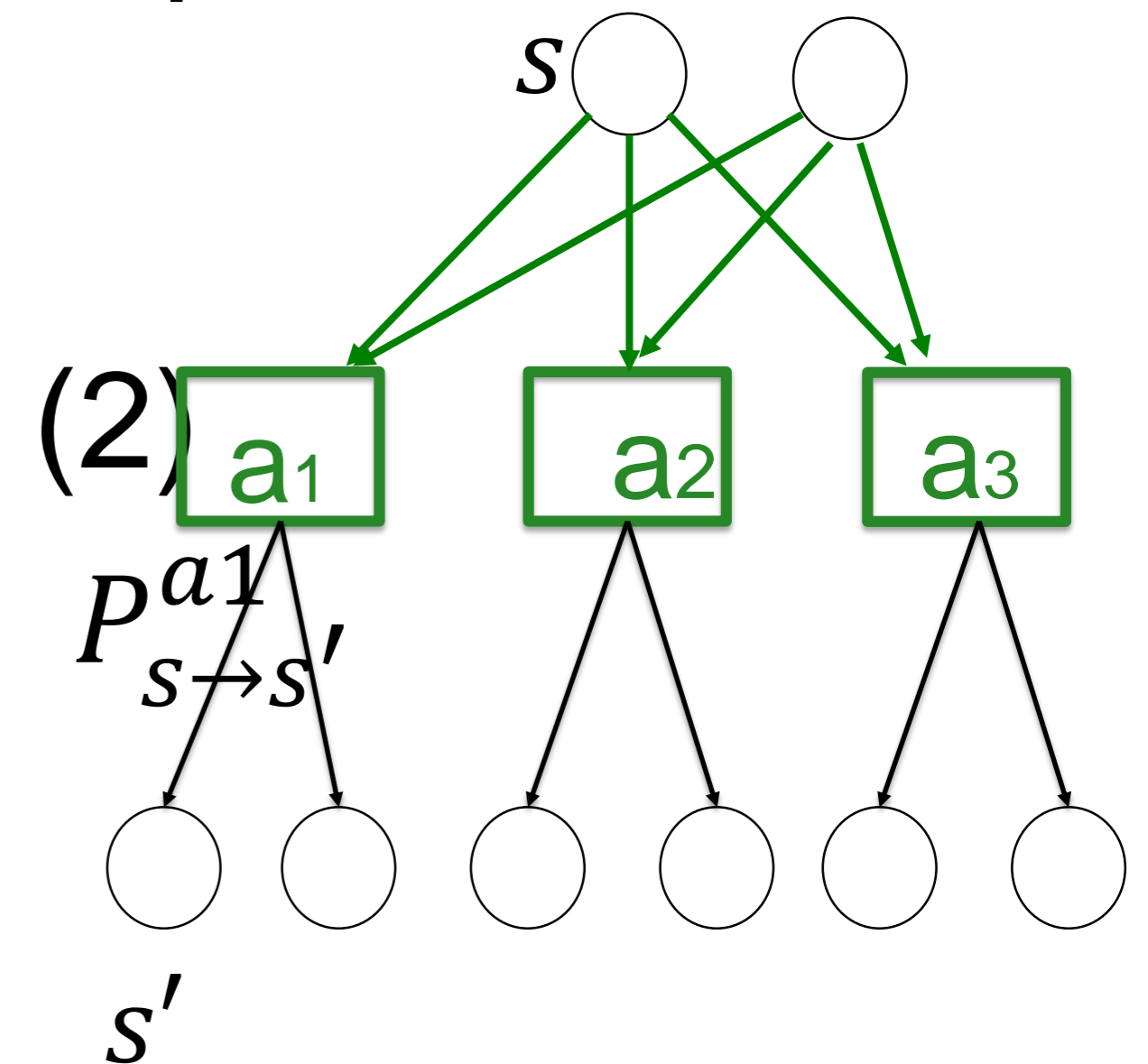
After taking action  $a$  in state  $s$ , we update with

$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)] \quad (1)$$

(i) If the expectation (over  $s'$ ,  $r$ ) of the update rule (1) **given (s,a)** vanishes, then  $\hat{Q}(s, a)$  has an expectation value,

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

(ii) If the learning rate  $\eta$  decreases, fluctuations around the **empirical mean**  $\langle \hat{Q}(s, a) \rangle$  decrease and the **empirical mean approaches**  $Q(s, a)$



Previous slide.

When evaluating the **expectation value given (s,a)**, the learning rate drops out since we set the left-hand-side to zero. The exact value of  $\eta$  is not relevant, as discussed in the theorem. Part (i) of the theorem states that the expectation value of  $\hat{Q}(s, a)$  is the correct Q-value. For a quick proof of part (i) see the video. On the blackboard a stronger statement was shown.

Convergence in expectation is equivalent to imagining that you start millions of trials with the same value  $\hat{Q}(s, a)$  without any intermediate update. So in that sense it is like a super-big 'batch' of examples.

**In practice, we do not have expectations but online updates with fluctuations.** It is important that the learning rate  $\eta$  is small at the end of learning so as to limit the amount of fluctuations. Part (ii) states that **online mean** for small learning rate also goes to the correct Q-value.

Indeed, since the equations are linear (for the bandit problem = 1-step horizon), the calculation of part (i) apply analogously to the long-term empirical temporal average (denoted by angular brackets)  $\langle \Delta \hat{Q}(s, a) \rangle = \eta \langle [r_t - \hat{Q}(s, a)] \rangle$

This equivalence based on linearity is not true for the multi-step horizon that we discuss later in this lecture.

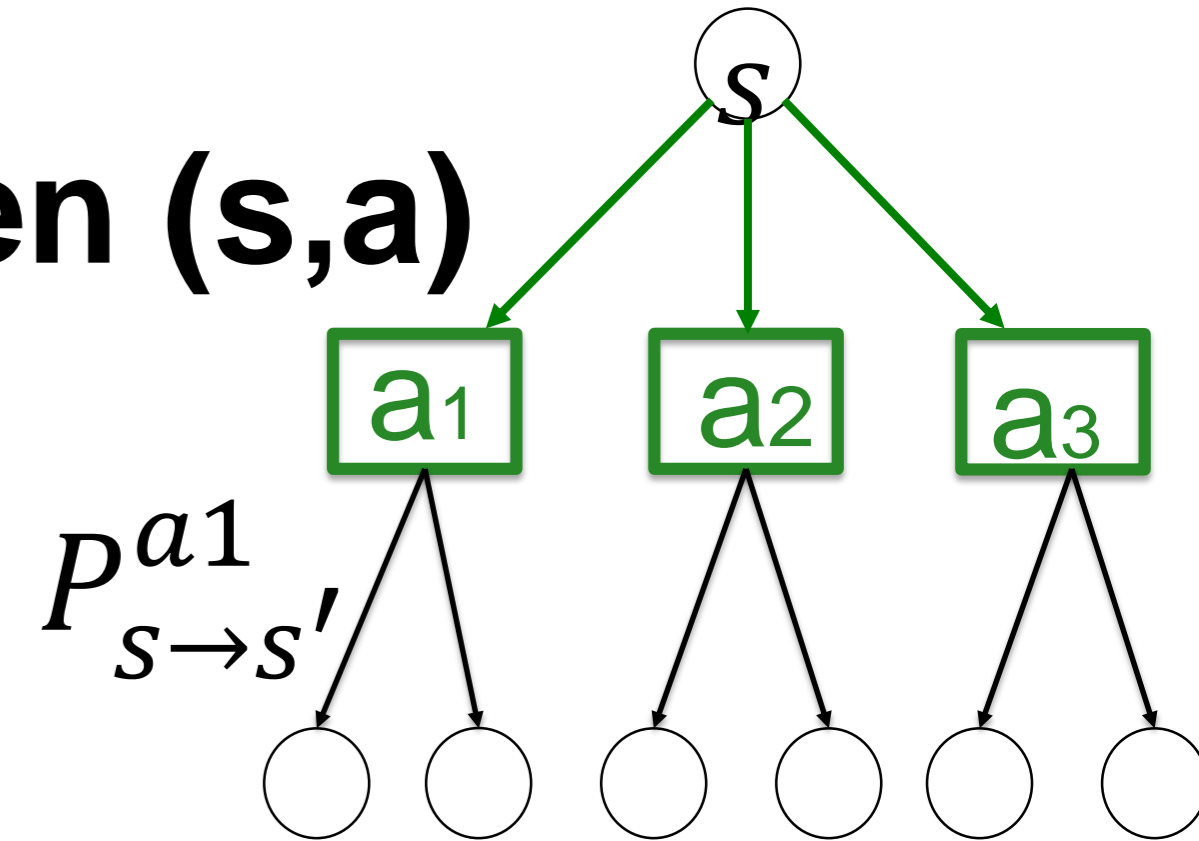
# Recall: Proof for strong version (blackboard last week)

After taking action  $a$  in state  $s$ , we update with

$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)] \quad (1)$$

(i) If the expectation of the update rule (1) given  $(s, a)$  vanishes, then  $\hat{Q}(s, a)$  has an expectation

$$\underline{E [\hat{Q}(s, a)] = \hat{Q}(s, a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a = Q(s, a)} \quad (2)$$



'I can even drop the expectation sign in the results'

If updates  $E(\Delta \hat{Q}(s, a) | s, a) = 0$   
then  $\hat{Q}(s, a)$  is the correct solution.

'  $\hat{Q}(s, a)$  takes  
role of  $\theta^{old}$  '

Your notes.

When I loosely say: 'converged in expectation' then I simply mean that the **expectation of the update rule vanishes.**

So the interpretation is:

if by chance we found a parameter  $\hat{Q}(s, a)$   
that fulfills

$$E(\Delta \hat{Q}(s, a) | \mathbf{s}, \mathbf{a}) = 0$$

then we know that  $\hat{Q}(s, a)$  is correct.

The theorem does not say how we would find it: we do not know whether it is stable under the online update algorithm.

# Part (i) of Theorem (strong version)

Expectation of online rule=0  $\rightarrow E(\Delta \hat{Q}(s, a) | s, a) = 0$

expectation of all possible futures with correct statistical weight

we always start in (s,a) while the system is frozen;  
 $\rightarrow \hat{Q}(s, a) = Q^{old}$  frozen

Update using **expectation over all possibilities that may occur in the next time step.**

**= Perspective similar to a batch mode:**

update only **after** (infinitely) many samples = trials that all start in (s,a) with the same frozen value  $Q^{old}$

$$E(\Delta \hat{Q}(s, a) | s, a) = \eta [E(r_t) - Q^{old}] = \lim_N \frac{1}{N} \sum_{k=1}^N \eta [r_t - Q^{old}]$$

Previous slide:

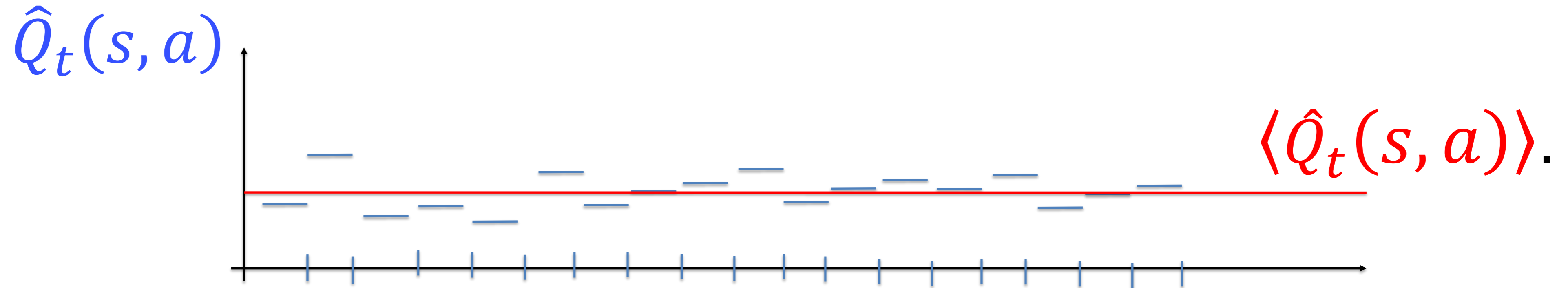
$\hat{Q}(s, a)$  denotes the current estimate of the Q-value. Claim: If Q no longer changes (in expectation) then it must be the correct Q-value.

There are different views on how to interpret the ‘expectation;:

- Formally from a mathematical point of view: average over all possible outcomes of the next time step given (s,a).
- In a simulation this would correspond to the following sampling procedure:  
You freeze the value of  $\hat{Q}(s, a)$  and run MANY times (N to infinity) a test with the state-action pair (s,a) as a starting condition. Then you evaluate the resulting ‘batch update’ averaged across all these examples. If the batch update with Millions of Examples is zero, that implies that you have converged to the correct value.

Part (ii) of Theorem:

We work with the **online update**  $\Delta \hat{Q}(s, a)$ . With finite learning rate, the value of  $\hat{Q}_t(s, a)$  fluctuates around a mean  $\langle \hat{Q}_t(s, a) \rangle$ .



**Claim:** Under the hypothesis  $\langle \Delta \hat{Q}(s, a) = 0 \rangle$ , the mean  $\langle \hat{Q}_t(s, a) \rangle$  is equal to the ‘correct’ Q-value.

Your notes. (Proof in the Blackboard notes)



Part (ii) of Theorem:

**Claim:** Under the hypothesis  $\langle \Delta \hat{Q}(s, a) = 0 \rangle$ , the temporal mean  $\langle \hat{Q}_t(s, a) \rangle$  with online updating is equal to the 'correct' Q-value  $Q(s, a)$ .

Proof:

$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}_t(s, a)] \quad (1)$$

$$\langle \Delta \hat{Q}(s, a) \rangle = \cancel{\eta} \langle r_t - \hat{Q}_t(s, a) \rangle = 0$$

$$0 = \langle r_t \rangle - \langle \hat{Q}_t(s, a) \rangle$$

$\langle r_t \rangle = Q(s, a)$  by definition of Q-values

Your notes. (Proof in the Blackboard notes)

# One-step horizon: summary

Q-value = expected reward for state-action pair

If Q-value is known, choice of action is simple

→ take action with highest Q-value

If Q-value not known:

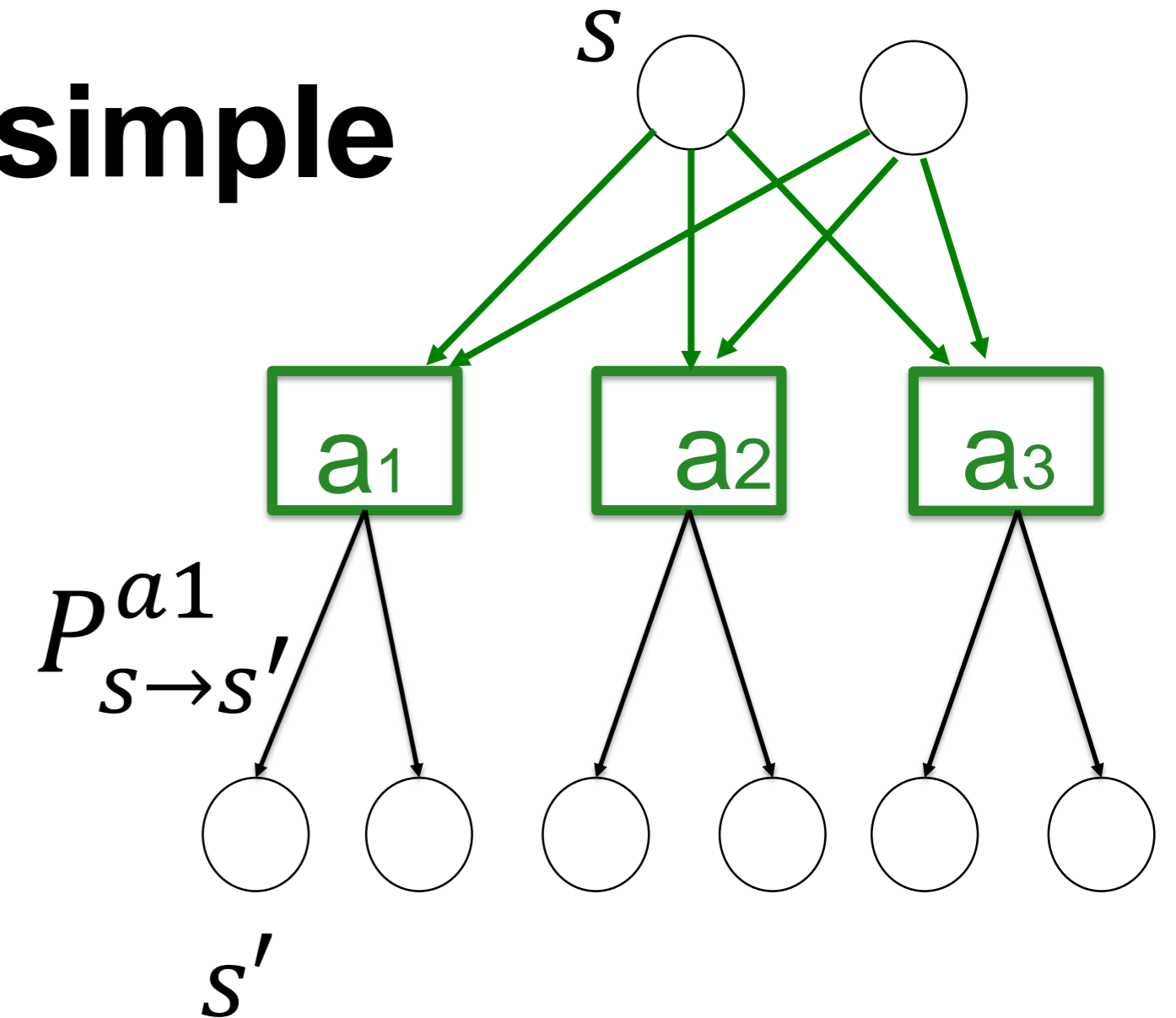
→ estimate  $\hat{Q}$  by trial and error

→ update with rule

$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)] \quad (1)$$

→ Let learning rate  $\eta$  decrease over time

- Iterative algorithm (1) always fluctuates for finite  $\eta$
- The expectation of the update step is 'good'



Previous slide.

Let us distinguish the ESTIMATE  $\hat{Q}(s, a)$  from the real Q-value  $Q(s, a)$

The update rule can be interpreted as follows:

if the actual reward is larger than (my estimate of) the expected reward, then I should increase (a little bit) my expectations.

The learning rate  $\eta$  :

In exercise 1, we found a rather specific scheme for how to reduce the learning rate over time. But many other schemes also work in practice. For example you keep  $\eta$  constant for a block of time, and then you decrease it for the next block.

**Note: I will often use the symbol  $\alpha$  instead of  $\eta$**

Both symbols indicate what is called the 'learning rate' in Deep Learning.

# Teaching monitoring – monitoring of understanding

The detour helped me to understand better and put the calculation from last week into perspective.

The detour was unnecessary.

The detour confused me even more

Previous slide.  
Teacher feedback

# Artificial Neural Networks: RL1

## Reinforcement Learning and SARSA

Wulfram Gerstner

EPFL, Lausanne, Switzerland

### Part 4: Exploration vs. Exploitation

- Examples of Reward-based Learning
- Elements of Reinforcement Learning
- One-step horizon (bandit problems)
- **Exploration vs. Exploitation**

Previous slide.

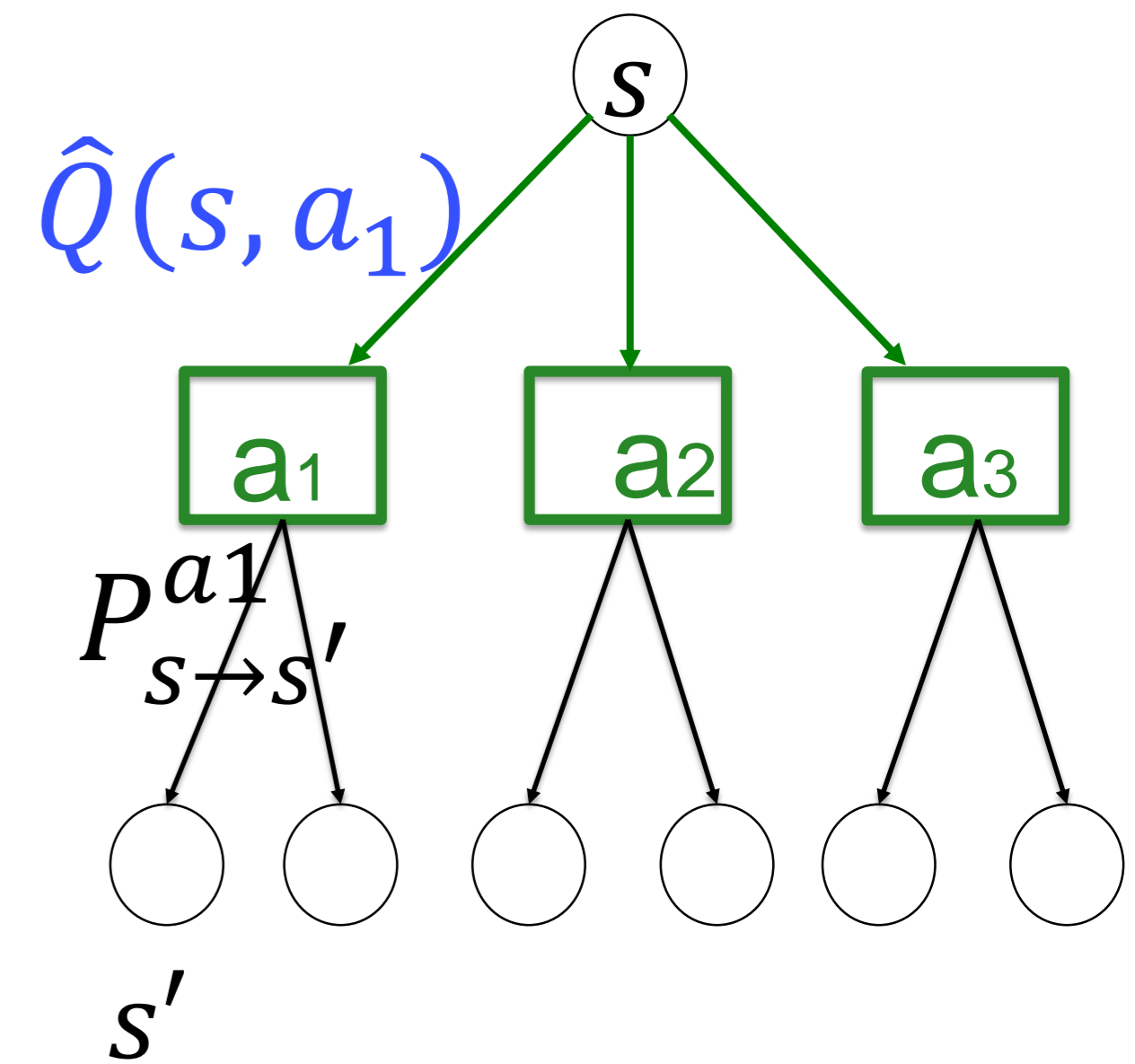
To estimate the Q-values you have to play all the different actions several times.  
However, if you know the Q-values you should only play the best action.



# Exploration – Exploitation dilemma

Ideal: take action with maximal  $Q(s, a)$

**Problem: correct Q values not known**  
(since reward probabilities and branching probabilities unknown)



## Exploration versus exploitation

Explore so as to estimate reward probabilities

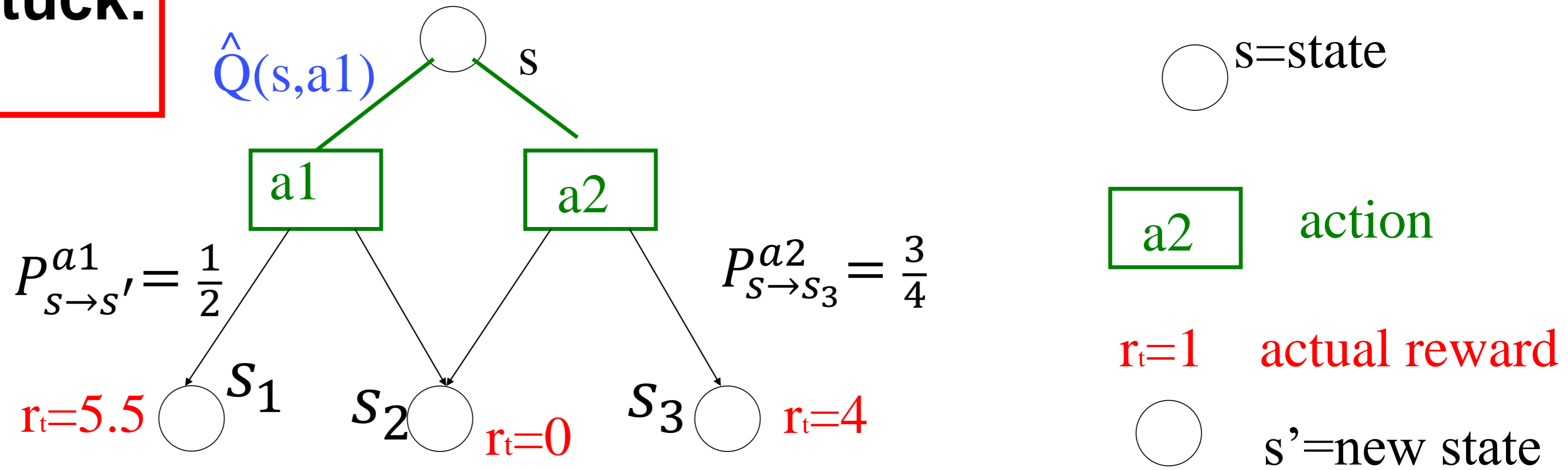
Take action which looks optimal, so as to maximize reward

Previous slide.

Since Q-values are not known, you are always in the situation of an exploration-exploitation dilemma.

Note: All estimates of Q will be empirical estimates. Here in this and the next slide I still write  $\hat{Q}$  for the empirical average. However, later, I simplify the notation and write for the empirical estimate  $Q(s,a)$  without the hat whenever the meaning is implicitly clear.

**greedy makes you stuck:  
Example**



Assume that you initialize all Q values with zero; set  $\eta = 0.2$  (constant)  
 update  $\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)]$

Trial 1: you choose action a1, you get  $r_t=5.5$

Trial 2: you choose action a2, you get  $r_t=4.0$



Trial 3 – 4: continue ‘greedy’:  $\rightarrow$  you continue with action 1

BUT: the expected reward  $Q(s, a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a$  is larger for action 2.

Given the outcomes of the first two trials, action a1 looks better.  
You can check that whatever the outcome in trial 3 (even for reward=0!), the estimated Q-value of action a1 is still higher than that of action a2!

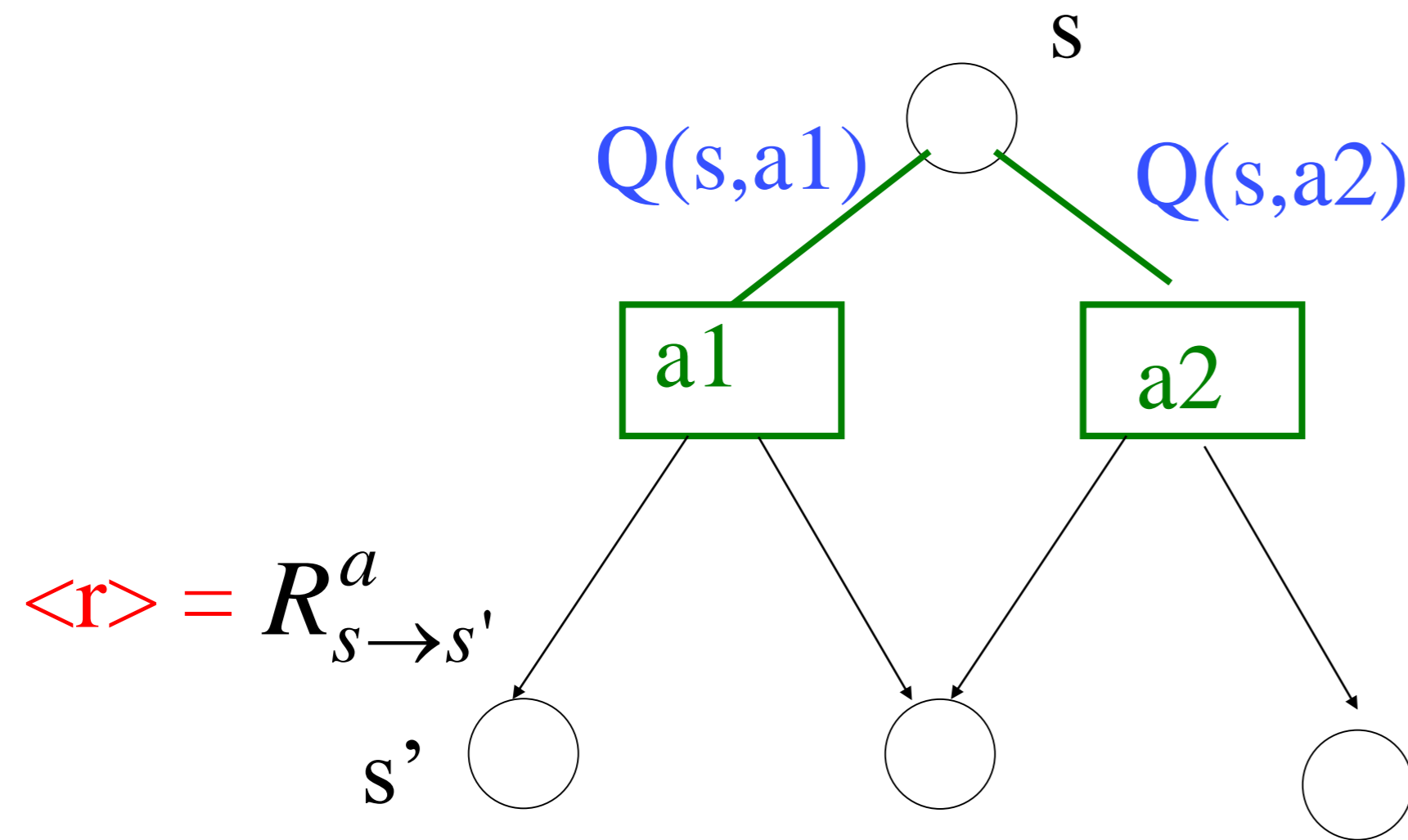
# Exploration and Exploitation

Problem: correct Q values not known

**greedy strategy:**

- take **action  $a^*$**  which looks best

$$Q(s, a^*) \geq Q(s, a_j) \quad \text{for all } j$$



**ATTENTION:**  
with 'greedy' you may get stuck with a sub-optimal strategy (see Exercise!)

Previous slide.

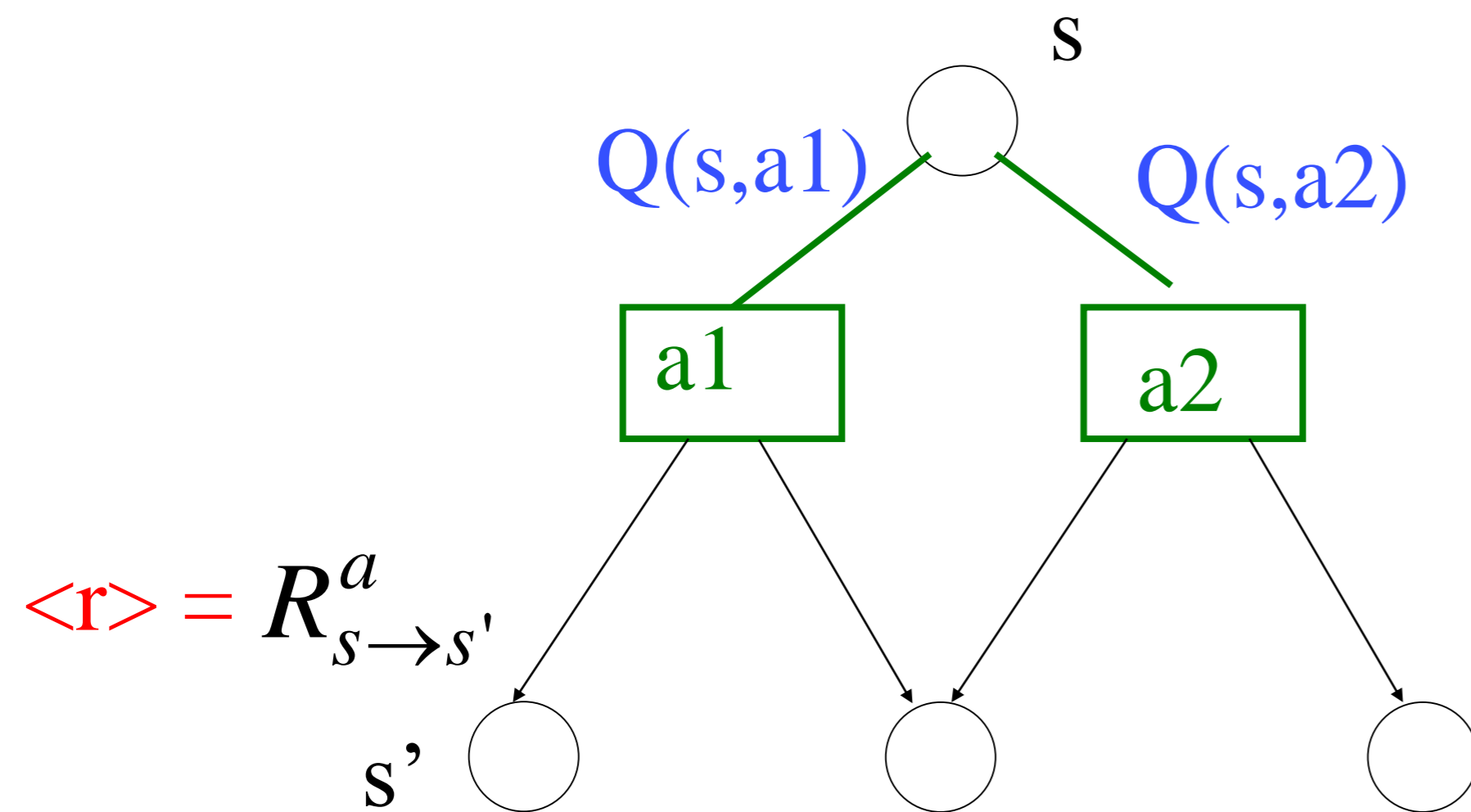
If you know the correct Q-values, the best choice would be to choose the action with maximal Q-value (called 'greedy' action). But since you don't know the Q-values it is risky to choose the greedy action because you may get stuck with a suboptimal choice.

In (almost all) applications of reinforcement learning we work with estimated Q-values.

Previously we used a hat to distinguish the ESTIMATED  $\hat{Q}(s, a)$  from the real Q-value  $Q(s, a)$ . However, in the following I will write the estimated Q-values without the hat. Nearly always Q means estimated Q.

# Exploration and Exploitation: practical approach

hats have been dropped!



Problem: correct Q values not known

**greedy strategy:**

- take **action  $a^*$**  which looks best

$$Q(s, a^*) \geq Q(s, a_j) \text{ for all } j$$

**$\epsilon$ -greedy strategy:**

- take **action  $a^*$**  which looks best

with prob  $P = 1 - \epsilon$

**Softmax strategy:** take **action  $a'$**

with prob

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

**Optimistic greedy:**

initialize with Q values that are too big

$$\Delta Q(s, a) = \eta [r_t - Q(s, a)]$$

Previous slide.

Softer versions of greedy allow you to choose occasionally an action which looks suboptimal, but which allows you to further explore the Q-values of other options.

Epsilon-greedy and softmax are examples following this idea.

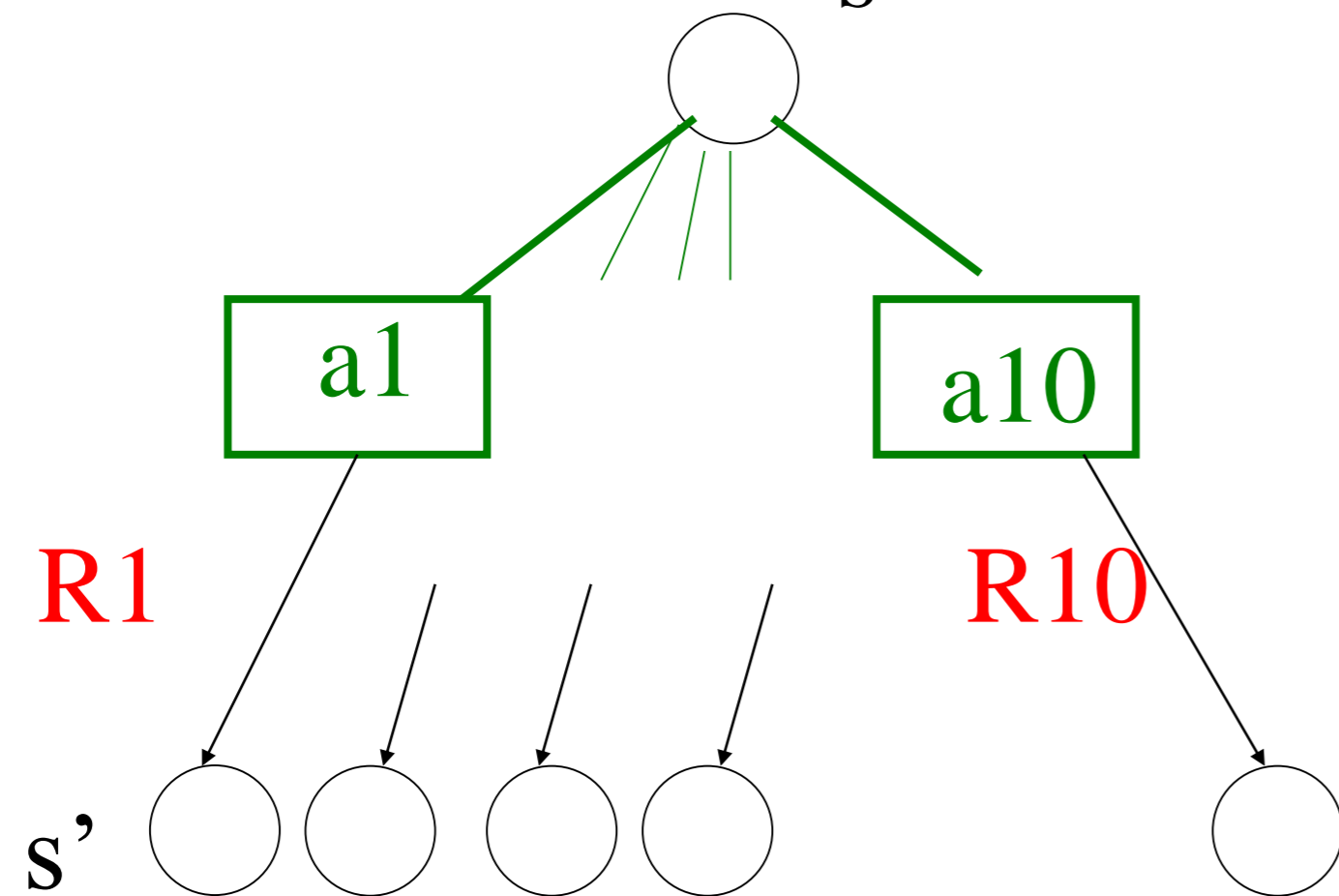
Note that 'softmax' is a function that one also encounters in multiclass tasks with 1-hot coding (see course of 'machine learning'; also later lecture on deep learning)

A radically different approach is optimistic greedy. If you initialize all Q-values at the same value, but clearly too high (compared to maximal reward that you can get in the scheme), then the Q-value of action  $a_1$  decreases initially each time you play  $a_1$ , which in turn favors other actions that you have not yet played.



# Exploration and Exploitation: practical approach

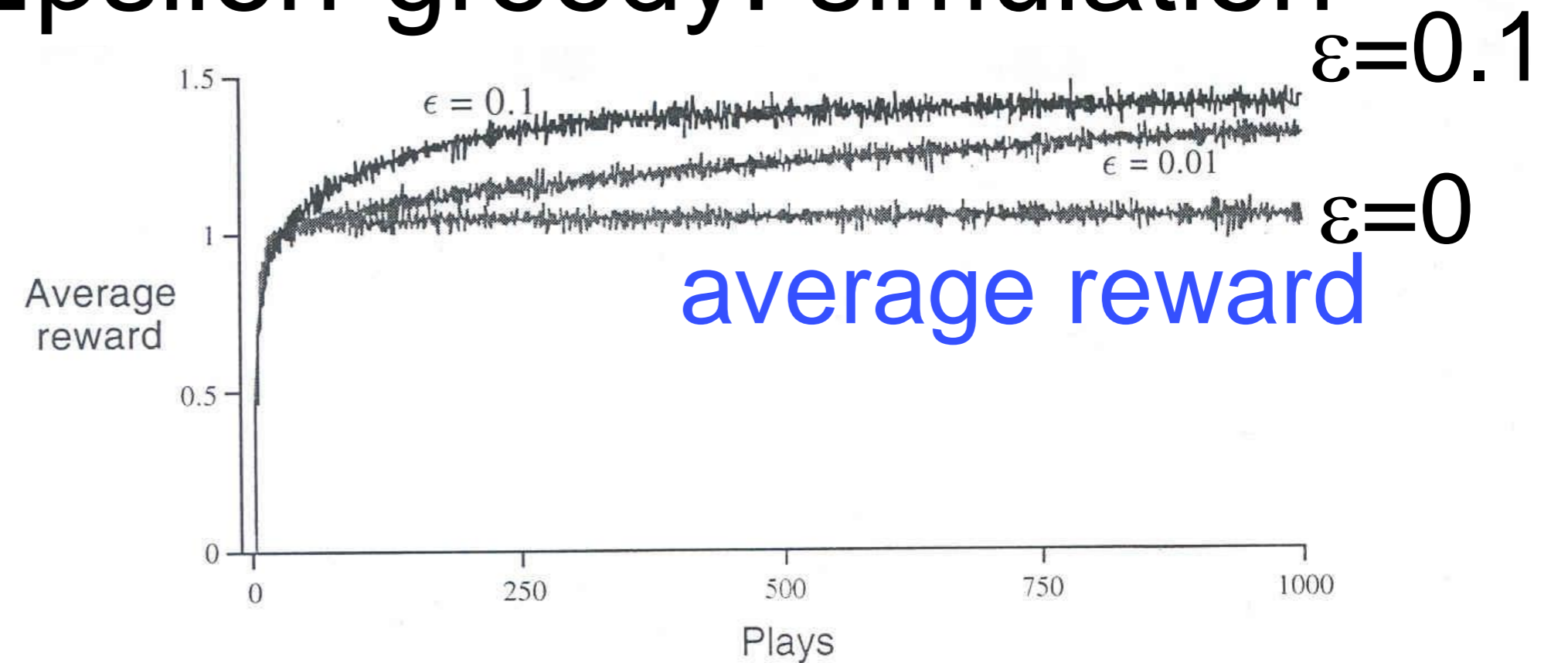
Example: 10-armed bandit with fluctuating reward



in each action, actual rewards fluctuate around a mean

$$R_k = R_{s \rightarrow s'}^{a_k}$$

Epsilon-greedy: simulation



Optimal action

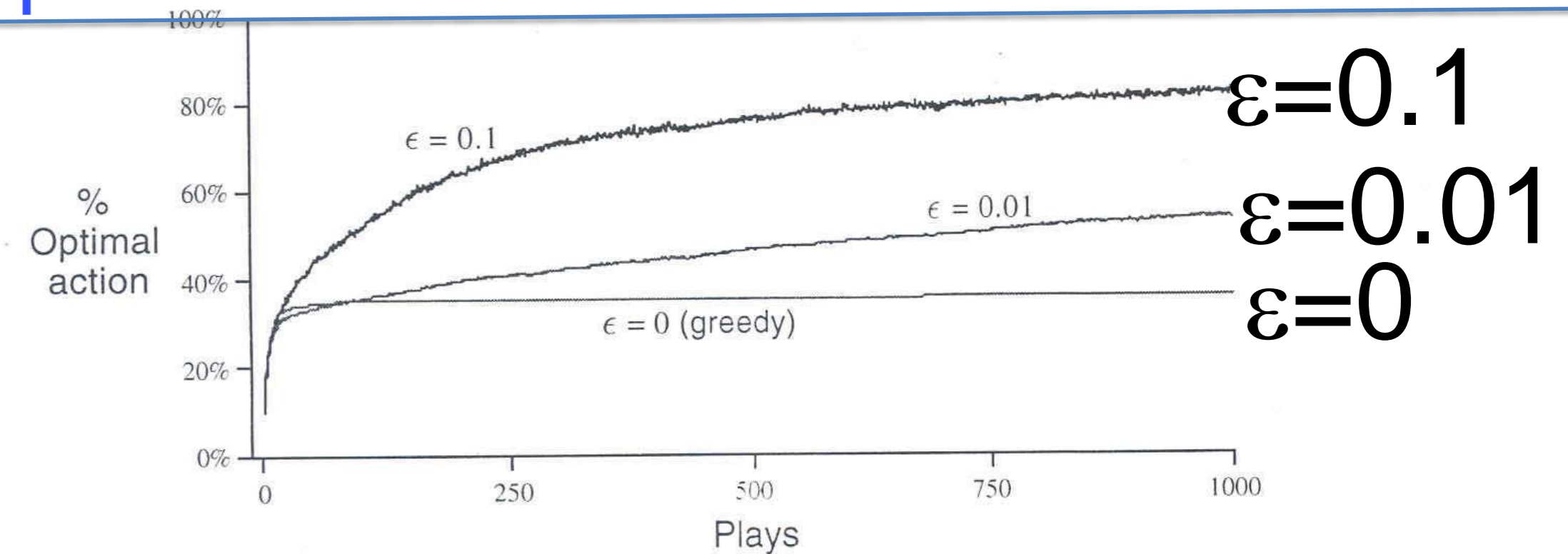


Figure 2.1 Average performance of  $\epsilon$ -greedy action-value methods on the 10-armed testbed. These data are averages over 2000 tasks. All methods used sample averages as their action-value estimates.

book: Sutton and Barto

Previous slide.

Computer simulation of a situation where actual rewards  $r$  fluctuate around the mean reward  $R$ . There are 10 different actions  $a_1, \dots, a_{10}$  each with a different mean reward  $R_1, \dots, R_{10}$ .

There exist two different ways to evaluate the performance.

Top: what is the average reward that you get by playing epsilon-greedy?

Bottom: what is the fraction of times that you play the optimal action, by playing epsilon-greedy.

Three different values of epsilon are used.

# Exploration and Exploitation: practical approach

Epsilon-greedy, combined with iterative update of Q-values

## A simple bandit algorithm

Initialize, for  $a = 1$  to  $k$ :

$$Q(a) \leftarrow 0$$

$$N(a) \leftarrow 0$$

Repeat forever:

$$A \leftarrow \begin{cases} \arg \max_a Q(a) & \text{with probability } 1 - \varepsilon \quad (\text{breaking ties randomly}) \\ \text{a random action} & \text{with probability } \varepsilon \end{cases}$$

$$R \leftarrow \text{bandit}(A)$$

$$N(A) \leftarrow N(A) + 1$$

$$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$$

learning rate

Sutton and Barto, ch. 2

Previous slide.

This is the style of pseudo-code that we will see a lot over the next few weeks. It is taken from the book of Sutton and Barto (MIT Press, 2018); Sutton and Barto made a pdf online available.

$Q(a)$  is the Q-value for action  $a$ . Since we have always the same starting state in which we have to make our choice of action, we can suppress the index of the state  $s$ .  $Q(a) = Q(s_{\text{start}}, a)$ .

$N(a)$  is a counter of how many times the agent has taken action  $a$ .

In this specific example the learning rate  $\eta$  is the inverse of the count  $N(a)$  (see earlier exercise); but in the more general setting we would remove the counter and just use some heuristic reduction scheme for  $\eta$ .

Note that in class we define  $(1-\epsilon)$  as the probability of taking the 'best' action corresponding to  $\text{argmax } Q$  and  $\epsilon$  is then distributed over the OTHER actions. Sutton and Barto distribute  $\epsilon$  over ALL actions, including the 'best'.

Thus for a total choice of 3 actions, Sutton and Barto have a probability of  $\epsilon/3$  for the other actions (and with the definition in class it would be  $\epsilon/2$ ).

# Quiz: Exploration – Exploitation dilemma

We use an iterative method and update Q-values with **eta=0.1**

With a greedy policy the agent uses the best possible action

Using an epsilon-greedy method with  $\epsilon = 0.1$  means that, even after convergence of Q-values, in about 10 percent of cases a suboptimal action is chosen.

If the rewards in the system are between 0 and 1 and Q-values are initialized with  $Q=2$ , then each action is played at least 5 times before exploitation starts.  
(exploitation starts when you no longer choose the wrong action)

Previous slide.

Here we define  $\epsilon$  as in class (1-epsilon) as the probability of taking the 'best' action corresponding to  $\text{argmax } Q$  and epsilon is then distributed over the OTHER actions.

# Quiz: Exploration – Exploitation with Softmax policy

Softmax policy: take **action  $a'$**  with prob  $P(a') = \frac{\exp[\beta Q(a')]}{\sum_a \exp[\beta Q(a)]}$

[ ] Suppose we have 3 possible actions  $a_1, a_2, a_3$  and use the softmax policy. Is the following claim true?

For  $Q(a_1) = 4, Q(a_2) = 1, Q(a_3) = 0$  the preference for action  $a_1$  is more pronounced

than for  $Q(a_1) = 34, Q(a_2) = 31, Q(a_3) = 30$ .

# Quiz: Exploration – Exploitation with Softmax policy

All Q values are initialized with the same value  $Q=0.1$

Rewards in the system are  $r = 0.5$  for action 1 (always)

and  $r=1.0$  for action 2 (always)

We use an iterative method and update Q-values with **eta=0.1**

1.  if we use softmax with **beta = 10**, then, after 100 steps, action 2 is chosen almost always
2.  if we use softmax with **beta = 0.1**, then, after 100 steps action 2 is taken about twice as often as action 1.

Softmax policy: take **action  $a'$**   
with prob  $P(a') = \frac{\exp[\beta Q(a')]}{\sum_a \exp[\beta Q(a)]}$



Your notes (Quiz not given in class).

Use that  $\exp(5)$  is a big number!

[yes], since  $\beta[Q(a_2) - Q(a_1)] = 5$

[no], with  $\beta = 0.1$ ,  $\exp(\beta Q) = 1 + \dots$   
→ both actions chosen with about the same prob.

Softmax policy: take **action  $a'$**   
with prob

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

# Exploration and Exploitation: Summary

- If we know the Q-values we can exploit our knowledge
- Exploitation = action which is best =  $\operatorname{argmax} Q(a)$
- But we never know the Q-values for sure
- We need to estimate the Q-values by playing the game
- Explore possibilities, transitions, outcomes, reward

**For complex problems, there is no perfect trade-off between exploration and exploitation**

Now Exercise 1\* - next lecture at 14h15

## Exercise: Exploration-Exploitation

Exercise from week 1

### Exercise 2. Greedy policy and the two-armed bandit

In the “2-armed bandit” problem, one has to choose one of 2 actions. Assume action  $a_1$  yields a reward of  $r = 1$  with probability  $p = 0.25$  and 0 otherwise. If you take action  $a_2$ , you will receive a reward of  $r = 0.4$  with probability  $p = 0.75$  and 0 otherwise. The “2-armed bandit” game is played several times and Q values are updated using the update rule  $\Delta Q(s, a) = \eta[r_t - Q(s, a)]$ .

- Assume that you initialize all Q values at zero. You first try both actions: in trial 1 you choose  $a_1$  and get  $r = 1$ ; in trial 2 you choose  $a_2$  and get  $r = 0.4$ . Update your Q values ( $\eta = 0.2$ ).
- In trials 3 to 5, you play greedy and always choose the action which looks best (i.e., has the

Your notes.

# Artificial Neural Networks: RL1

## Reinforcement Learning and SARSA

Wulfram Gerstner

EPFL, Lausanne, Switzerland

### Part 5: Bellman Equation

- Examples of Reward-based Learning
- Elements of Reinforcement Learning
- One-step Horizon (Bandit Problems)
- Exploration vs. Exploitation
- **Bellman Equation**

Previous slide.

So far our Q-values were limited to situations with a 1-step horizon. Now we will get more general.

# Teaching monitoring – monitoring of understanding

[ ] today, up to here, at least 60% of material was new to me.

[ ] up to here, I have the feeling that I have been able to follow (at least) 80% of the lecture.

Remark.

1. As an EPFL teacher, my aim is to teach such that 80 percent of the students in the classroom are able to follow at least 80 percent of the lecture.

Sometimes I succeed, sometimes I don't.

2. Since most courses in the master are optional, overlap is sometimes unavoidable. However, my aim as an EPFL teacher is to still present in each lecture 60 percent of material that is new to 60 percent of the students.

If part of a lecture is novel for less than 20 percent of the students, I am happy to remove that material from the in-class presentation (it could still stay in the notes).

So tell me, if you see a large overlap with other classes, and then I can take a poll in class to find out how many people are concerned.



# Multistep horizon

**Policy**  $\pi(s, a)$

probability to choose  
action  $a$  in state  $s$

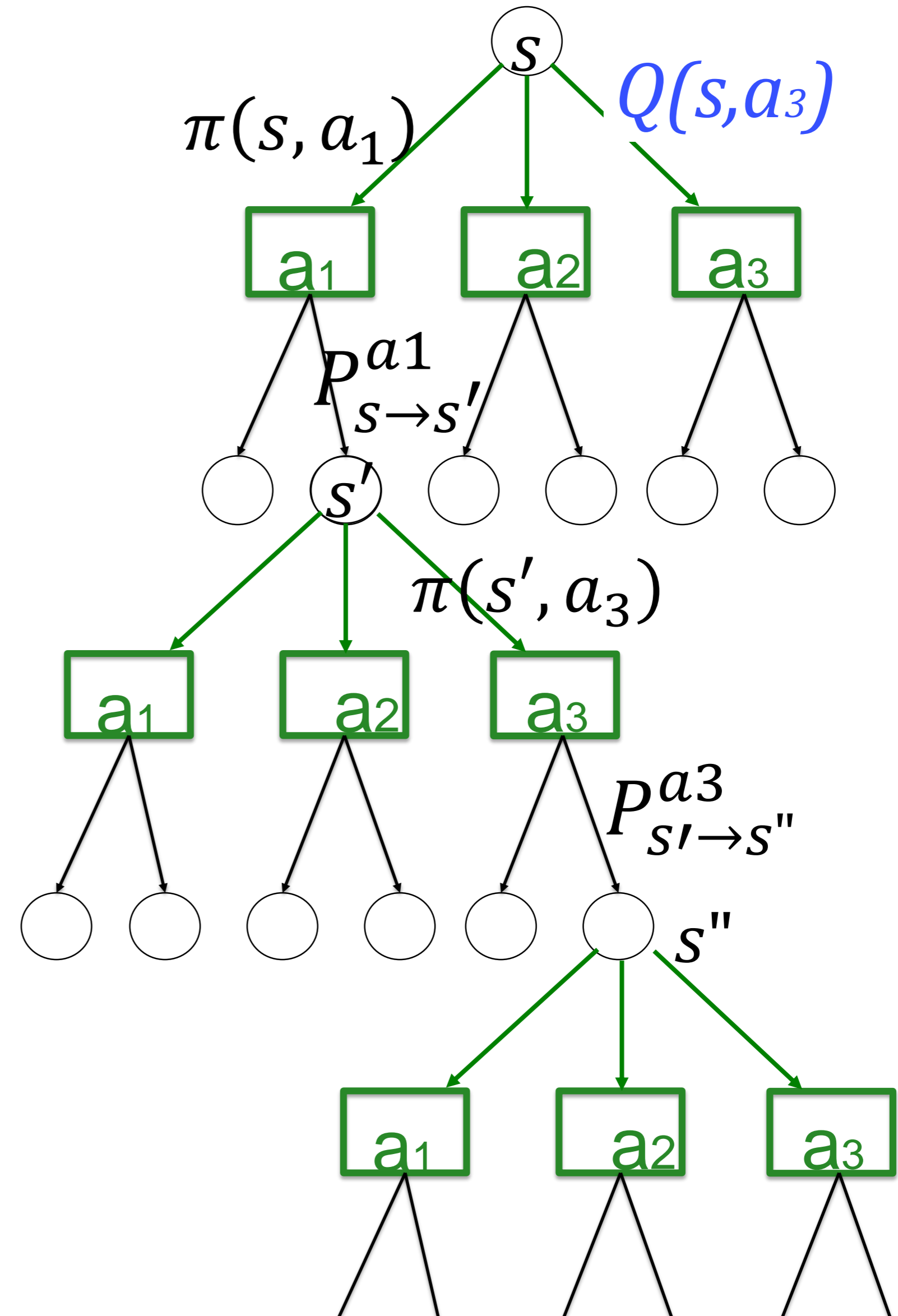
$$1 = \sum_{a'} \pi(s, a')$$

**Examples of policy:**

- epsilon-greedy
- softmax

**Stochasticity**  $P_{s \rightarrow s'}^a$

probability to end in state  $s'$   
taking action  $a$  in state  $s$



Previous slide.

After a first action that leads to state  $s'$  starting from state  $s$ , the agent can now take a second action starting from  $s'$ .

Note that there are two different types of branching ratio:

$\pi(s, a_1)$  describes the probability that the agent uses action  $a_1$  when it is in state  $s$  – based on the agent's policy (such as epsilon-greedy)

$P_{s \rightarrow s'}^{a_1}$  describes as before the probability that the agent arrives in state  $s'$  given that it chooses action  $a_1$  in state  $s$ .

As before we are interested in the expected reward. The Q value  $Q(s,a)$  describes the total accumulated reward the agent can get starting in state  $s$  with action  $a$ .

Next slide: rewards that are  $n$  steps away are discounted with a factor  $\gamma^n$

# Total expected (discounted) reward

Starting in state  $s$  with action  $a$

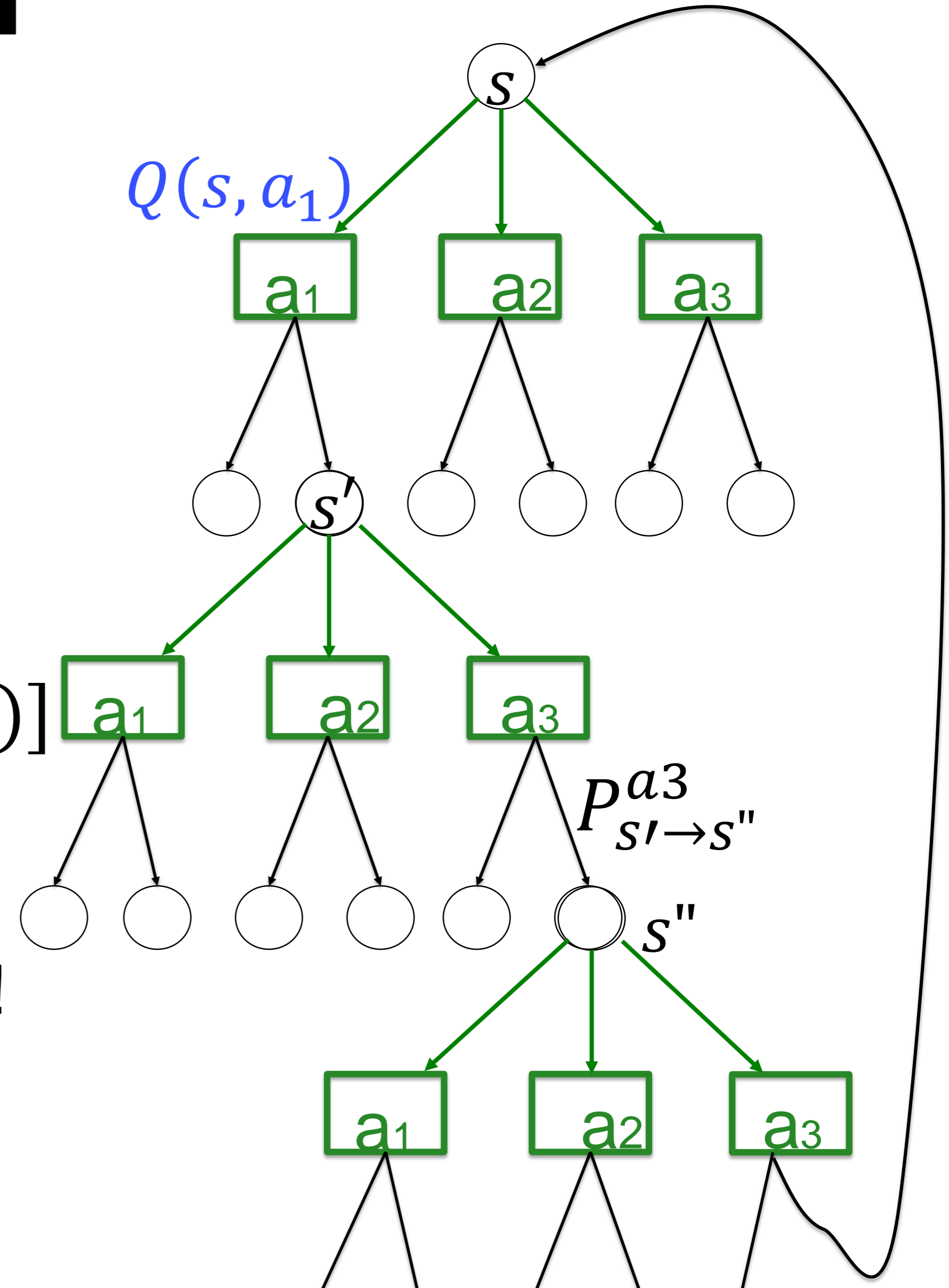
$$Q(s,a) =$$

$$= \langle r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots \rangle$$

$$= E[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots | s, a]$$

**Discount factor:  $\gamma < 1$**

- important for recurrent state transition graphs!
- avoids blow-up of summation
- gives less weight to reward in **far future**



Previous slide.

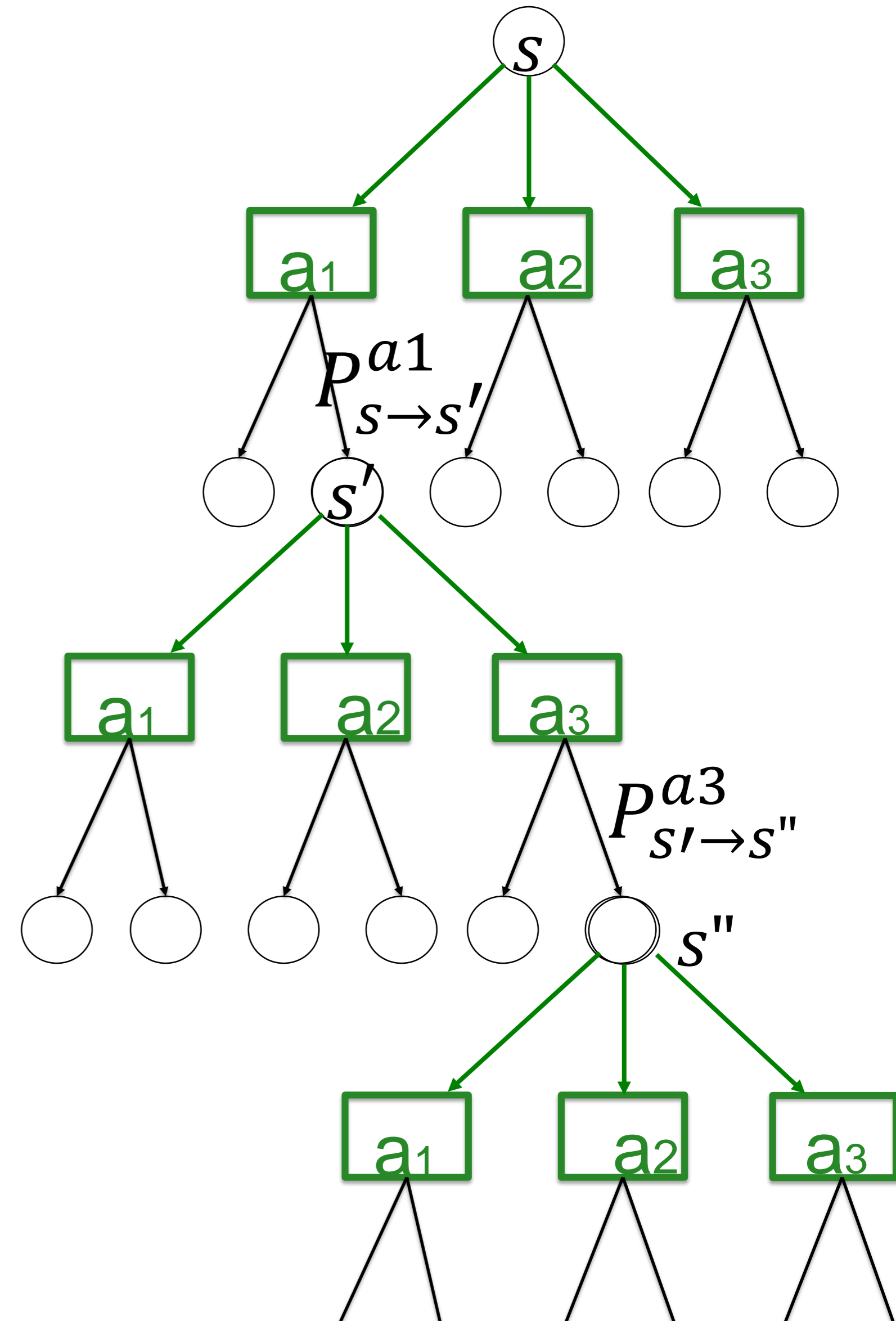
Angular brackets denote expectation.

Red-font lower-case  $r$  indicates the reward collected over multiple time steps in one episode, starting in state  $s$  with action  $a$ .

Expectation means that we have to take the average over all possible future paths giving each path its correct probabilistic weight.

# Bellman equation

Blackboard4:  
Bellman eq.



Space for calculations.

# Bellman equation with policy $\pi$

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

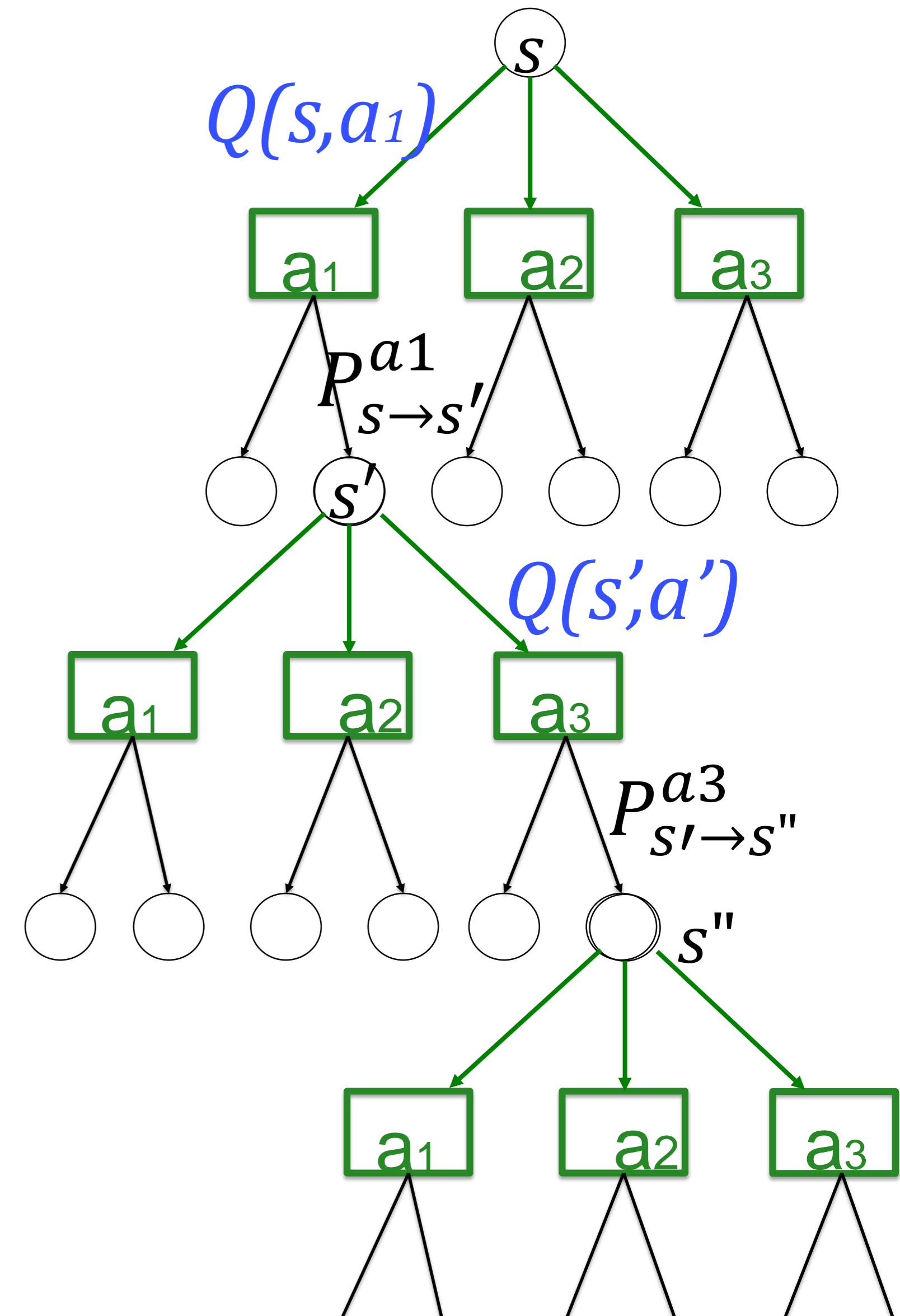
Bellman equation =  
value consistency of  
neighboring states

## Remark:

Sometimes Bellman equation is written

for greedy policy:  $\pi(s, a) = \delta_{a, a^*}$

with action  $a^* = \operatorname{argmax}_{a'} Q(s, a')$



Previous slide.

The Bellman equation relates the Q-value for state  $s$  and action  $a$  with the Q-values of the neighboring states.

Neighboring means reachable in a single step.

Note that the two different types of branching ratio both enter the equation.

Bottom: in the case of a greedy policy, the Bellman equation simplifies



# Bellman equation (for optimal actions)

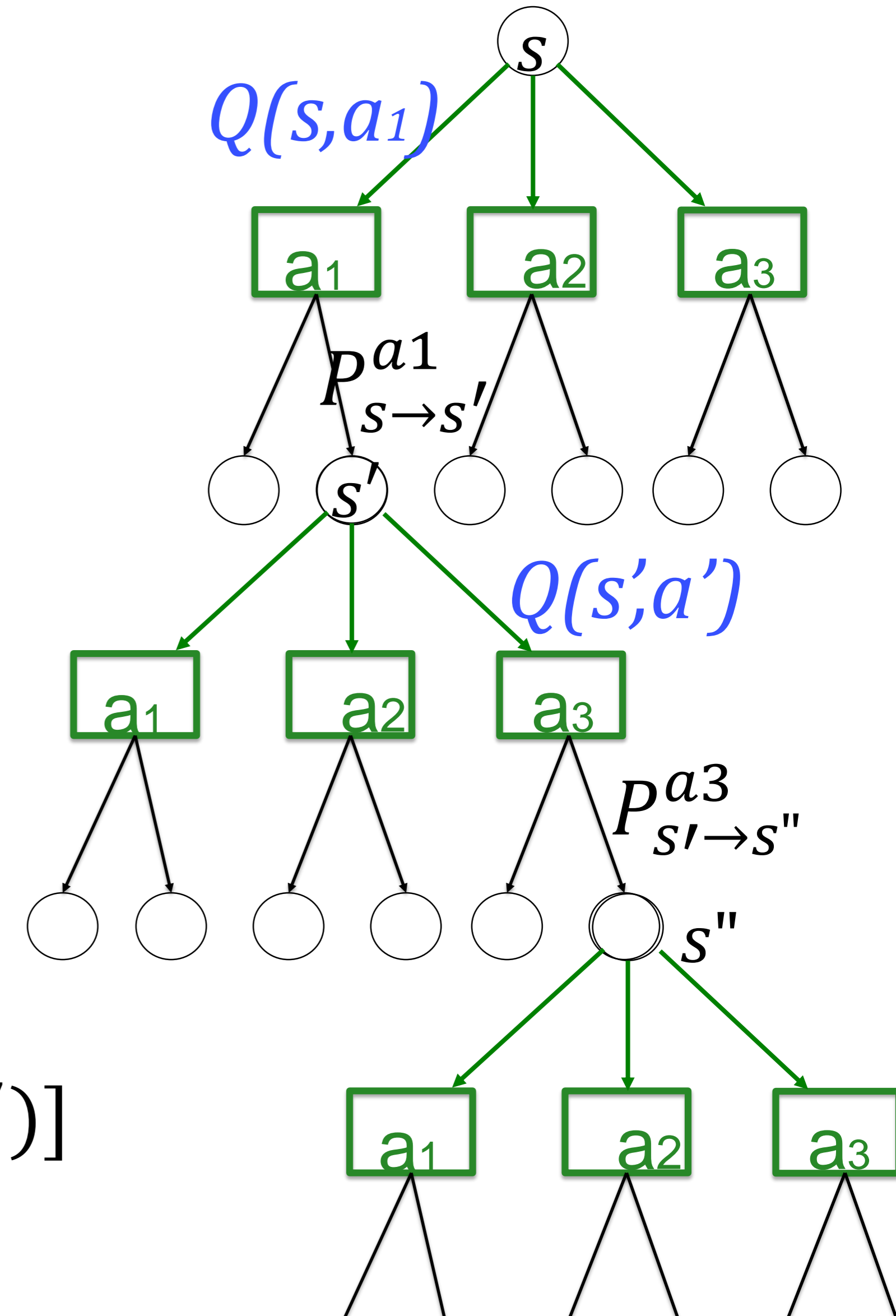
$$Q(s, a) = \sum_{s'} P_{s \rightarrow s'}^a \left[ R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi(s', a') Q(s', a') \right]$$

for greedy policy:

$$\pi(s, a) = \delta_{a, a^*}$$

with action  $a^* = \operatorname{argmax}_{a'} Q(s, a')$

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



Previous slide.

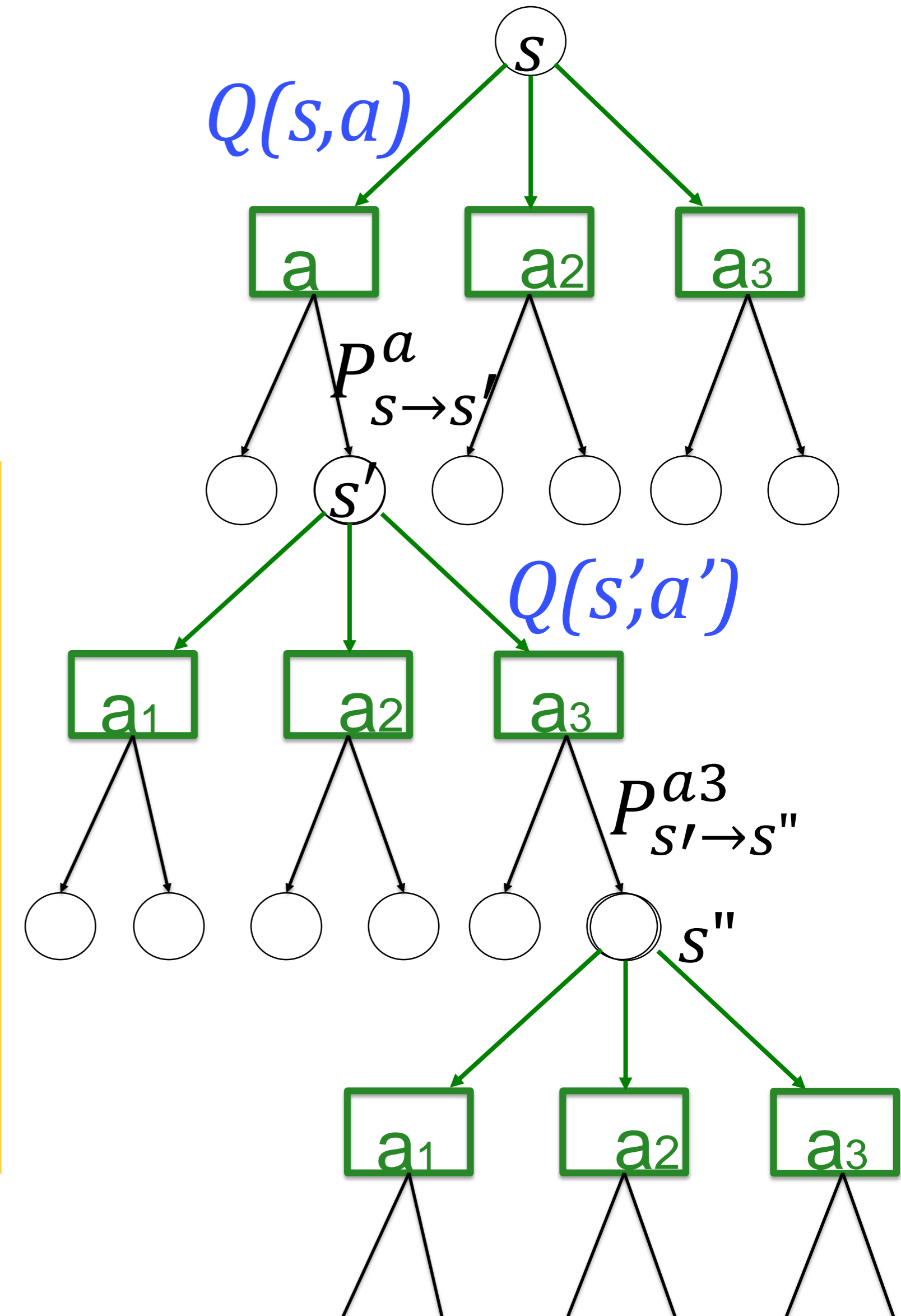
For a greedy policy, the sum over actions disappears from the Bellman equation and is replaced by the max-sign.

# Quiz: Bellman equation with policy $\pi$

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

[ ] The Bellman equation is linear in the variables  $Q(s', a')$

[ ] The set of variables  $Q(s', a')$  that solve the Bellman equation is unique and does not depend on the policy



Your comments.

# Artificial Neural Networks: RL1

## Reinforcement Learning and SARSA

Wulfram Gerstner

EPFL, Lausanne, Switzerland

### Part 6: SARSA Algorithm

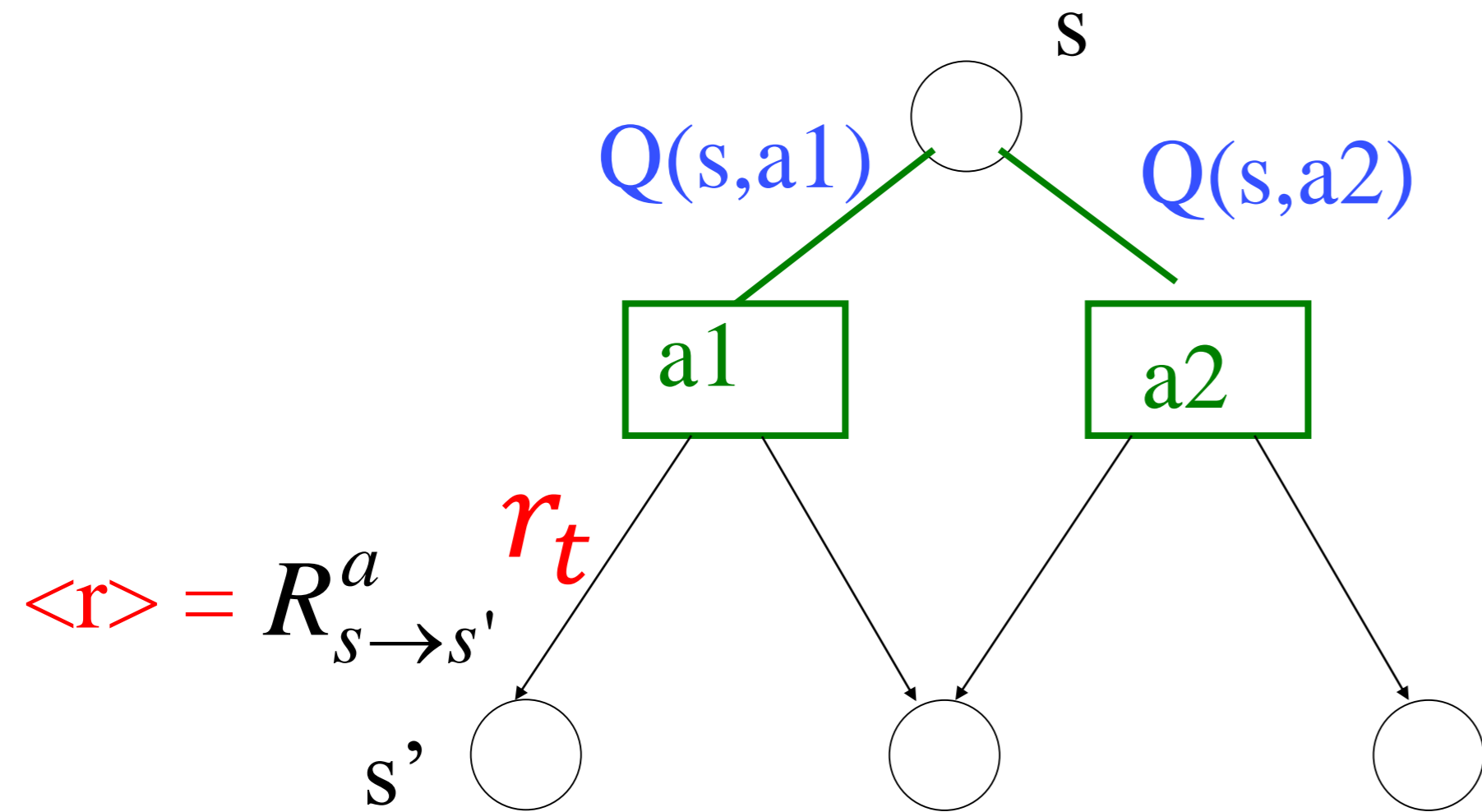
- Examples of Reward-based Learning
- Elements of Reinforcement Learning
- One-step Horizon (Bandit Problems)
- Exploration vs. Exploitation
- Bellman Equation
- **SARSA Algorithm**

Previous slide.

We not turn to the first practical algorithm, called SARSA. This is an algorithm that is widely used in the field of reinforcement learning.

# Review: Iterative update of Q-values

**Problem:** Q-values not given



**Solution:** iterative update

$$\Delta Q(s, a) = \eta [r_t - Q(s, a)]$$

while playing with policy  $\pi(s, a)$

Previous slide.

Reminder: for the 1-step horizon scenario we found that we could calculate the Q-values iteratively.

We increase the Q-value by a small amount (with learning rate  $0 < \eta \ll 1$ ) if the reward observed at time  $t$  is larger than our current estimate of  $Q$ .

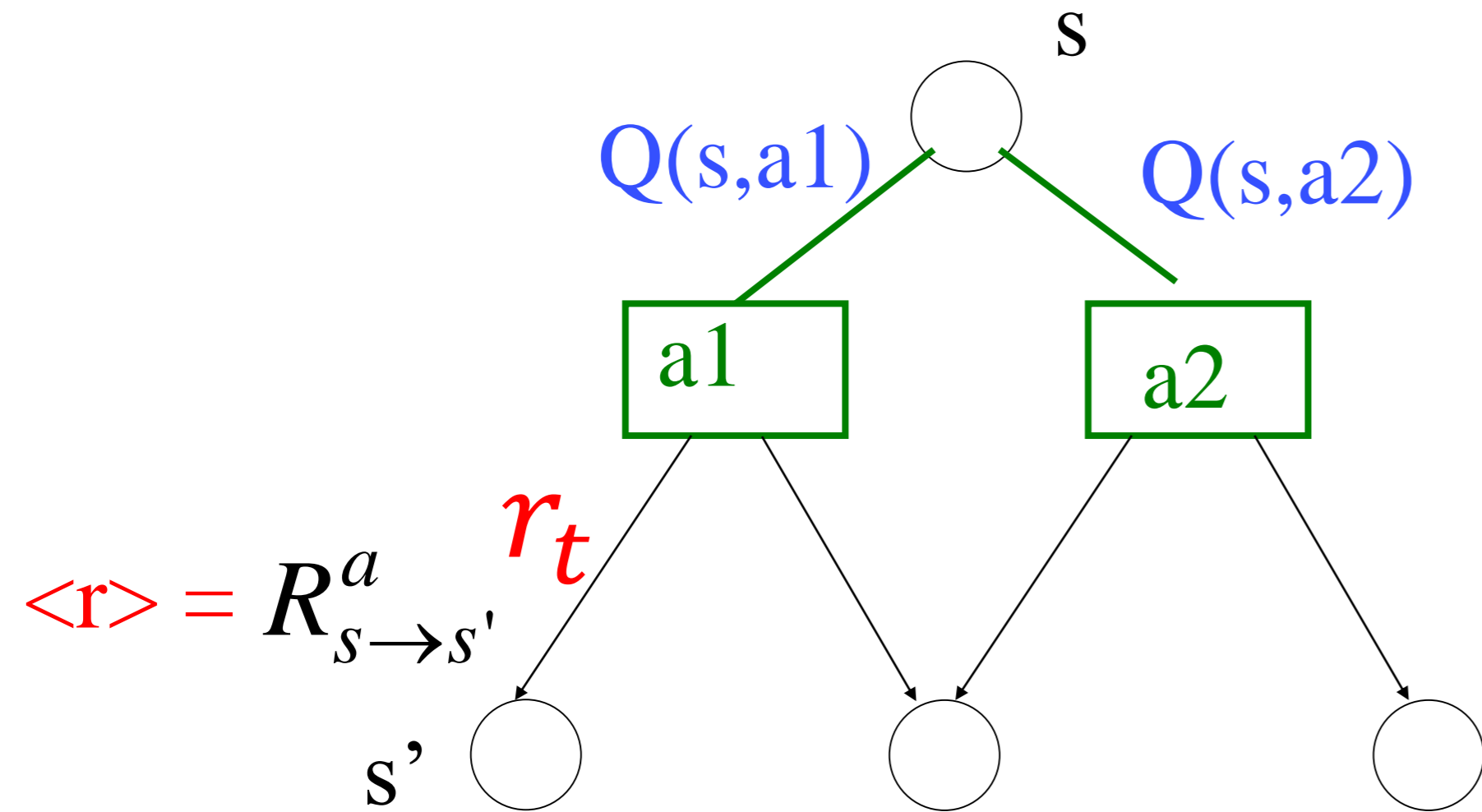
And we decrease the Q-value by a small amount if the reward observed at time  $t$  is smaller than our current estimate of  $Q$ .

Iterative updates with one data point at a time are also called 'online algorithms'. Thus our update rule is an online algorithm for the estimation of Q-values.



# Iterative update of Q-values for multistep environments

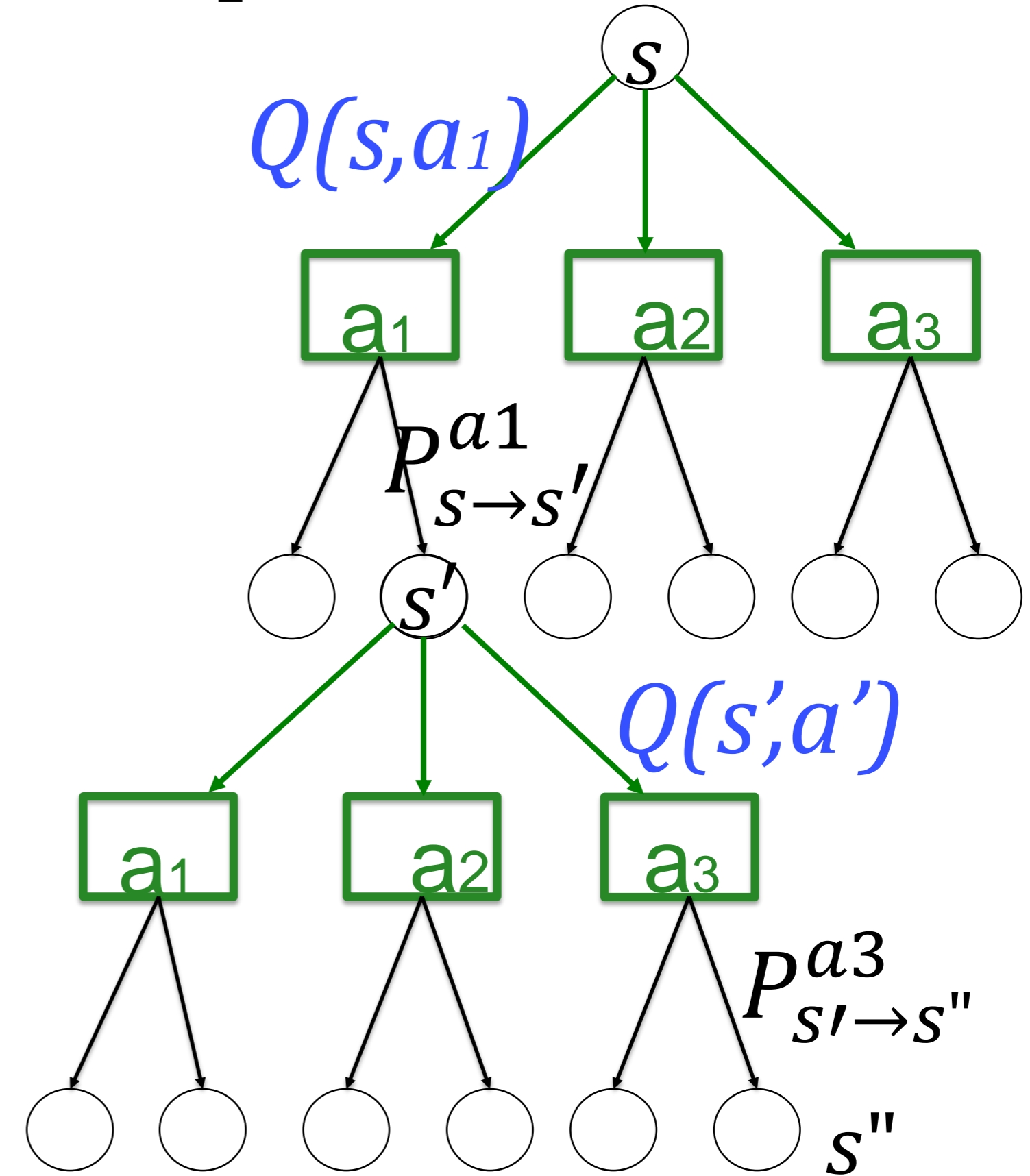
**Problem:** Q-values not given



**Solution:** iterative update

$$\Delta \hat{Q}(s, a) = \eta [r_t - \hat{Q}(s, a)]$$

while playing with policy  $\pi(s, a)$

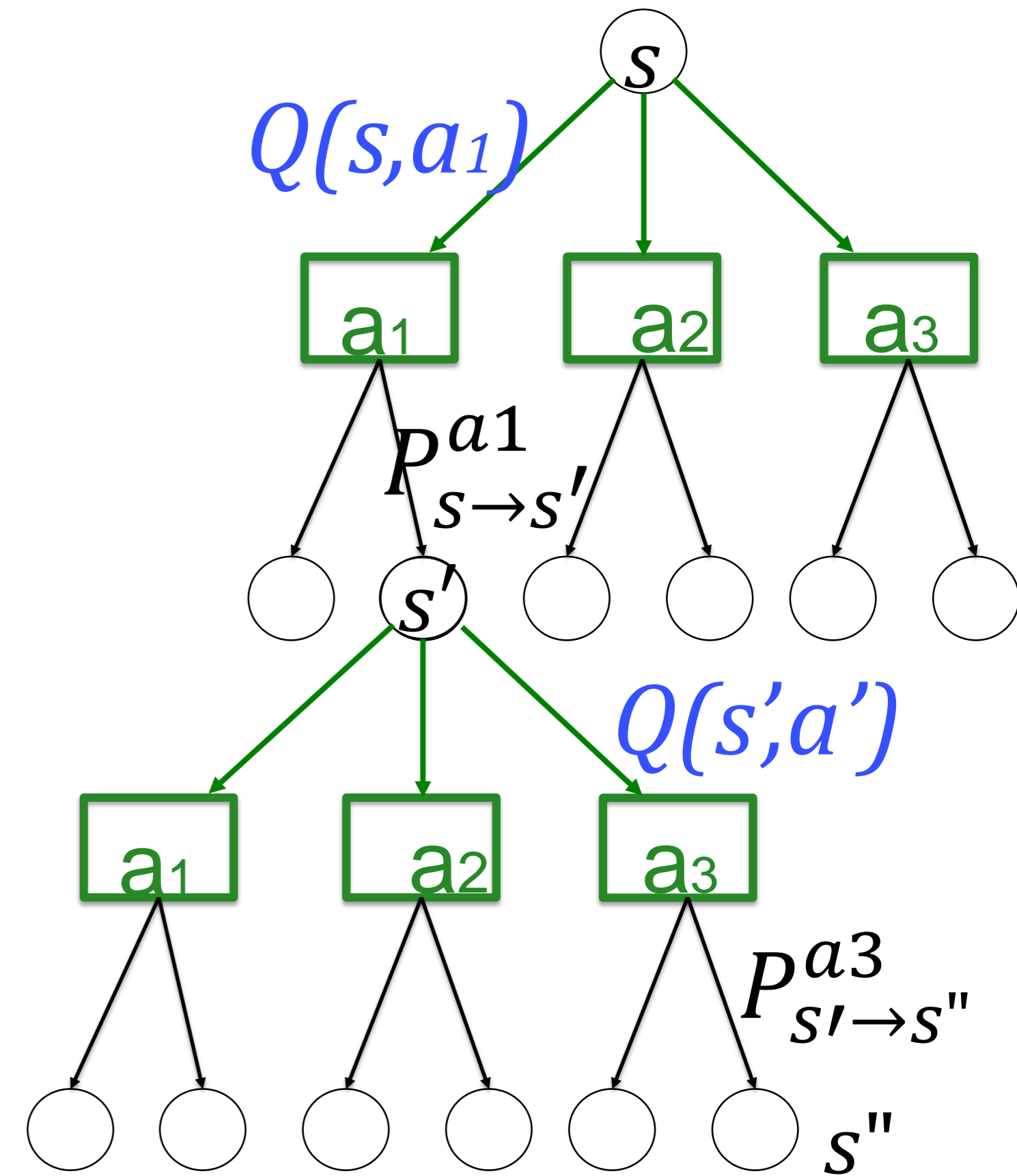


$$\Delta \hat{Q}(s, a) = ?$$

Previous slide.

The question now is: can we have a similar iterative update scheme also for the multi-step horizon?

# Blackboard5: SARSA update



Your notes.

# Iterative update of Q-values for multistep environments

## Bellman equation:

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

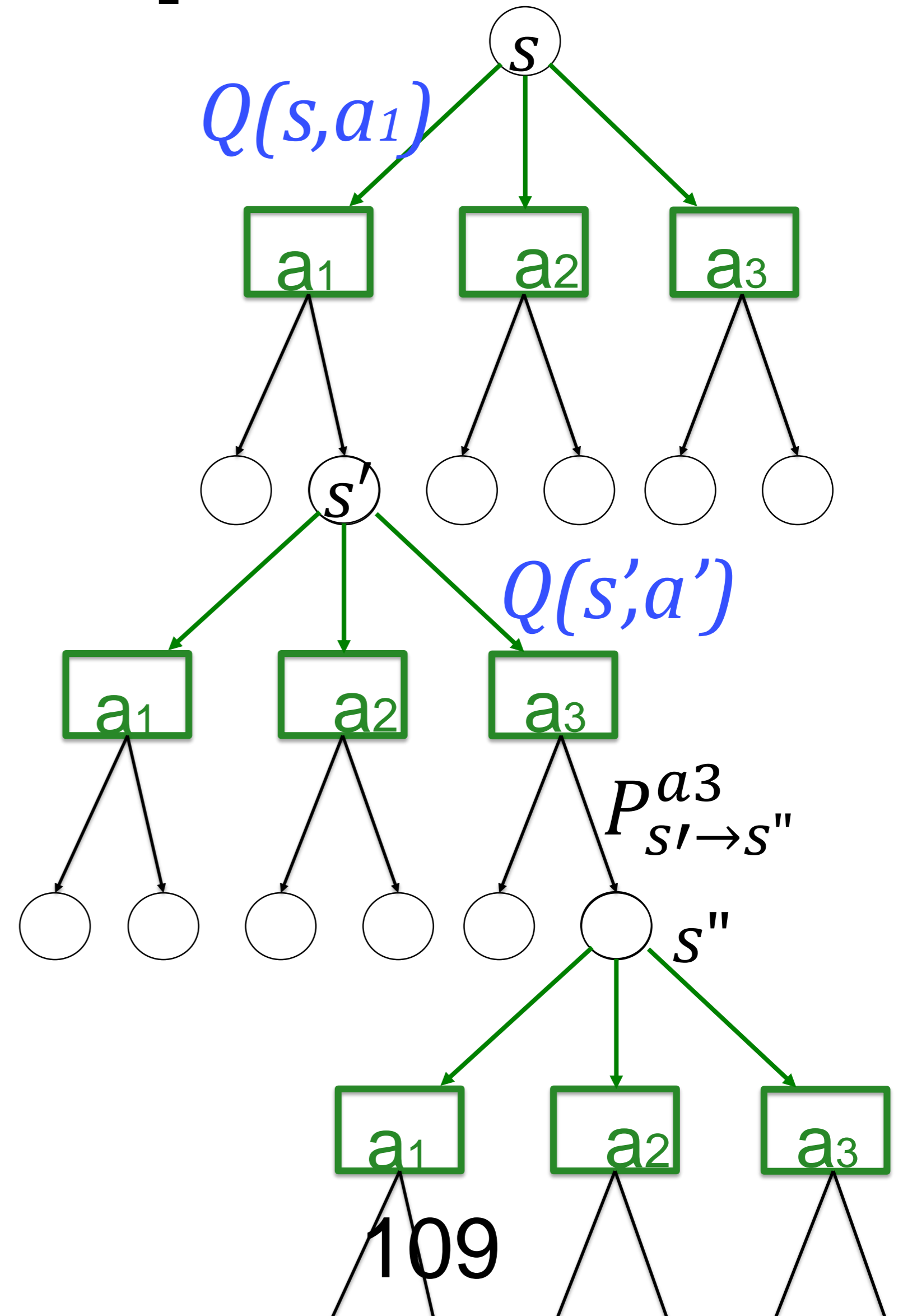
## Problem:

- Q-values not given
- branching probabilities not given
- reward probabilities not given

## Solution: iterative update

$$\Delta \hat{Q}(s, a) = \eta [r_t + \gamma \hat{Q}(s', a') - \hat{Q}(s, a)]$$

while playing with policy  $\pi(s, a)$



Previous slide.

Even for the case of the multi-step horizon, we can estimate the Q-values by an iterative update:

The Q-values  $Q(s,a)$  is increased by a small amount if the sum of (reward observed at time  $t$  plus discounted Q-value in the next step) is larger than our current estimate of  $Q(s,a)$ .

This iterative update gives rise to an online algorithm.

NOTE: in the following we always work with empirical estimates, and drop the 'hat' of the variable  $Q$ .

# SARSA vs. Bellman equation

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

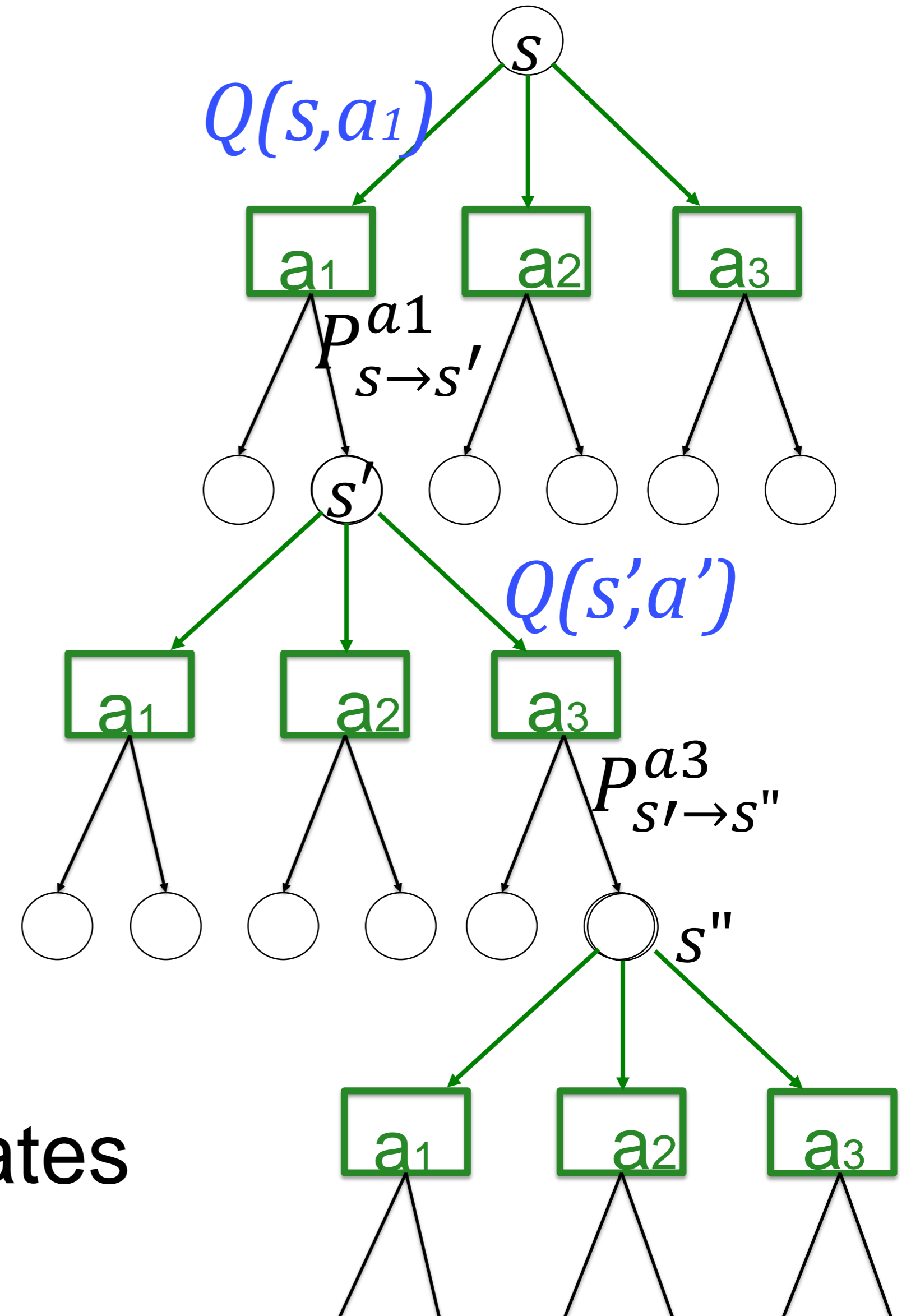
## Bellman equation

= consistency of Q-values across neighboring states

## SARSA update rule

$$\Delta Q(s, a) = \eta [r_t + \gamma Q(s', a') - Q(s, a)]$$

= make Q-values of neighboring states more consistent



Previous slide.

The Bellman equation summarizes the consistency condition:

The (average) rewards must explain the difference between  $Q(s,a)$  and  $\gamma Q(s',a')$  averaged over all  $s'$  and  $a'$ .

Or equivalently:

$Q(s,a)$  must be explained by the (average) reward in the next step and the discounted Q-value in the next state.

The iterative update formula implies that  $Q(s,a)$  needs to be adapted so the current reward explains the difference between  $Q(s,a)$  and  $Q(s',a')$ .

An equivalent form of writing the update is:

$$\Delta Q(s,a) = \eta [r - (Q(s,a) - \gamma Q(s',a'))]$$

This form highlights the similarity to the case of the update rule in the case of the one-step horizon.



# SARSA algorithm

Initialise Q values

Start from initial state  $s$

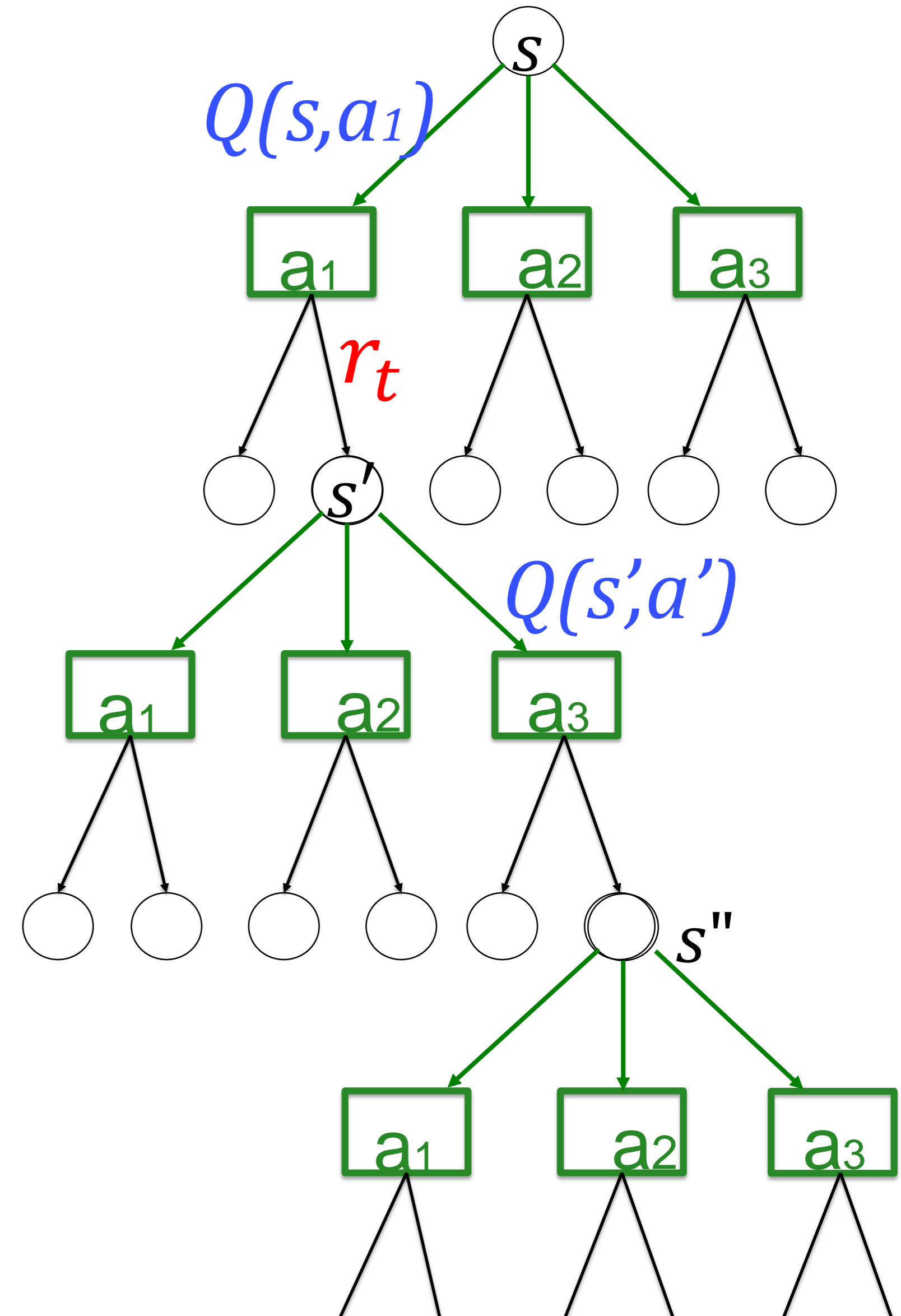
- 1) being in state  $s$   
and having chosen action  $a$   
[according to policy  $\pi(s, a)$ ]
- 2) Observe reward  $r$   
and next state  $s'$
- 3) Choose action  $a'$  in state  $s'$   
[according to policy  $\pi(s, a)$ ]
- 4) Update with SARSA update rule

$$\Delta Q(s, a) = [r_t + \gamma Q(s', a') - Q(s, a)]$$

5) set:  $s \leftarrow s'$ ;  $a \leftarrow a'$

6) Goto 1)

Stop when all Q-values have converged



Previous slide.

The update rule gives immediately rise to an online algorithm. You play the game. While you run through one of the episodes you observe the state  $s$ , choose action  $a$ , observe reward  $r$ , observe next state  $s'$  and choose next action  $a'$ . At this point in time (and not earlier) you have all the information to update the Q-value  $Q(s,a)$ .

The name SARSA comes from this sequence state-action-reward-state-action.

# SARSA algorithm.

$$\Delta Q(s, a) = \eta [r_t + \gamma Q(s', a') - Q(s, a)]$$

We have initialized SARSA and played for  $n > 2$  steps.

Is the following true for the next steps?

in SARSA, updates are applied after each move.

in SARSA, the agent updates the Q-value  $Q(s^{(t)}, a^{(t)})$  related to the **current** state  $s^{(t)}$

in SARSA, the agent updates the Q-value  $Q(s^{(t-1)}, a^{(t-1)})$  related to the **previous** state, once it has chosen  $a^{(t)}$

in SARSA, the agent moves in the environment using the policy  $\pi(s, a)$

SARSA is an online algorithm

Previous slide.

# Exercise at 15h15: SARSA in 1-dim environment

- Update of Q values in SARSA (with  $\gamma = 1$ )

$$\Delta Q(s, a) = [r_t + \gamma Q(s', a') - Q(s, a)]$$

- policy for action choice:

Pick most often action

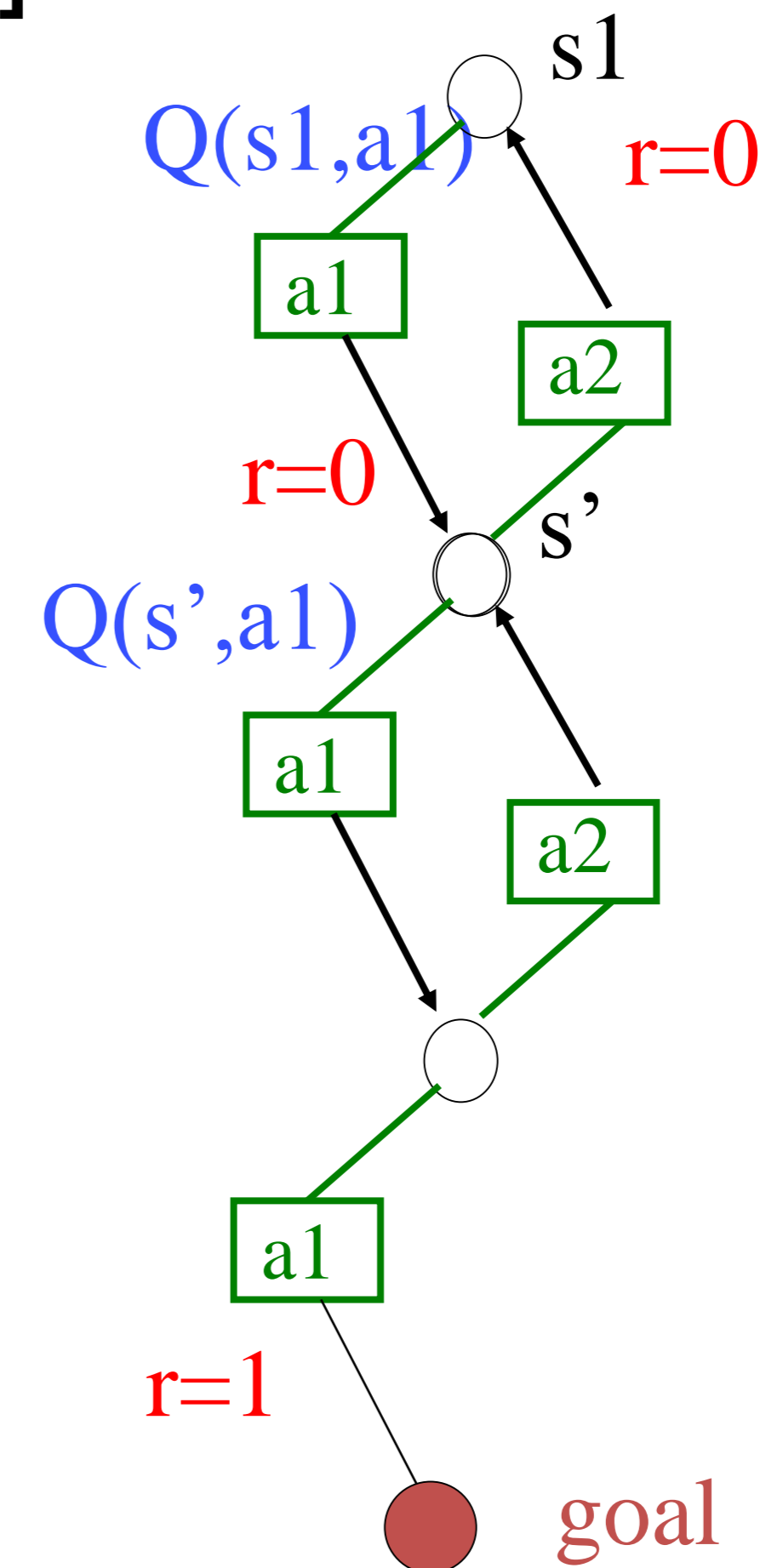
$$a_t^* = \arg \max_a Q_a(s, a)$$

Break ties stochastically

- States form linear sequence
- Reward only at goal.

- Initialise Q values at 0. Start at top state s1.

- Q values after 2 complete episodes?



# Exercise 4 (at 15h15)

## Exercise 4: SARSA for Linear Track.

### Exercise 4. SARSA algorithm

In the lecture, we introduced the SARSA (state-action-reward-state-action) algorithm, which (for discount factor  $\gamma = 1$ ) is defined by the update rule

$$\Delta Q(s, a) = \eta [r - (Q(s, a) - Q(s', a'))] , \quad (1)$$

where  $s'$  and  $a'$  are the state and action subsequent to  $s$  and  $a$ . In this exercise, we apply a greedy policy, i.e., at each time step, the action chosen is the one with maximal expected reward, i.e.,

$$a_t^* = \arg \max_a Q_a(s, a) . \quad (2)$$

If the available actions have the same Q-value, we take both actions with probability 0.5.

Consider a rat navigating in a 1-armed maze (=linear track). The rat is initially placed at the upper end of the maze (state  $s$ ), with a food reward at the other end. This can be modeled as a one-dimensional sequence of states with a unique reward ( $r = 1$ ) as the goal is reached. For each state, the possible actions are going up or going down (Fig. 2). When the goal is reached, the trial is over, and the rat is picked up by the experimentalist and placed back in the initial position  $s$  and the exploration starts again.

- Suppose we discretize the linear track by 6 states,  $s_1, \dots, s_6$ . Initialize all the Q-values at zero. How do the Q-values develop as the rat walks down the maze in the first trial?
- Calculate the Q-values after 3 complete trials. How many Q-value values are non-zero? How many trial do we need so that information about the reward has arrived in the state just 'below' the starting state?
- What happens to the learning speed if the number of states increases from 6 to 12? How many Q-values are non-zero after 3 trials? How many trial do we need so that information about the reward has arrived in the state just 'below' the starting state?

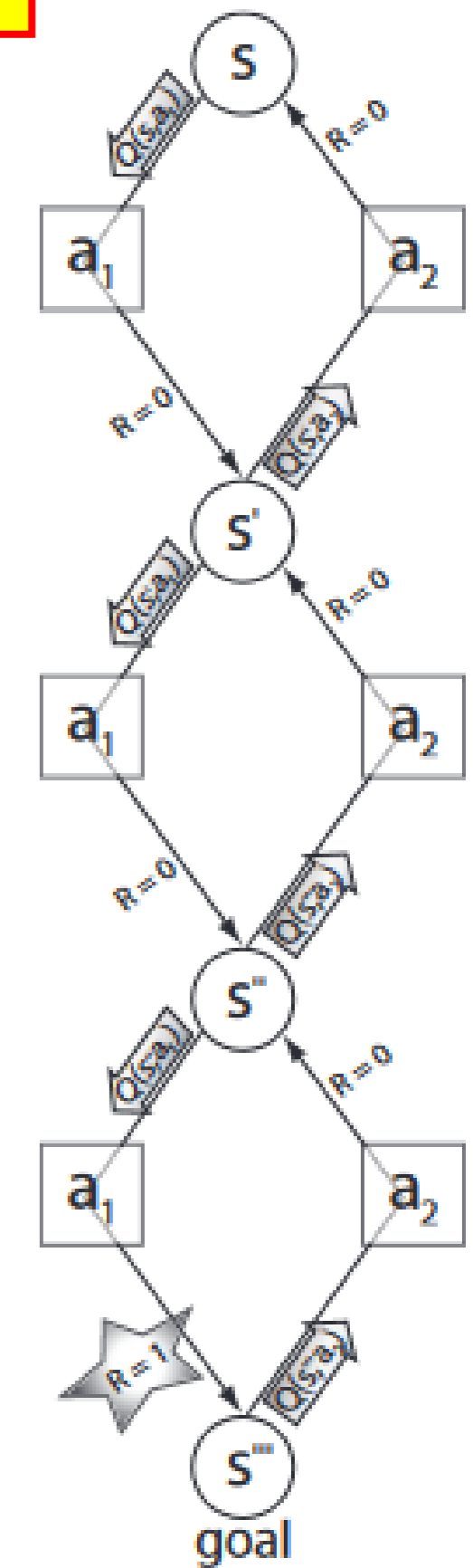


Figure 2: A linear maze.

# Variant A: SARSA is consistent w. Bellman equation

We prove the following:

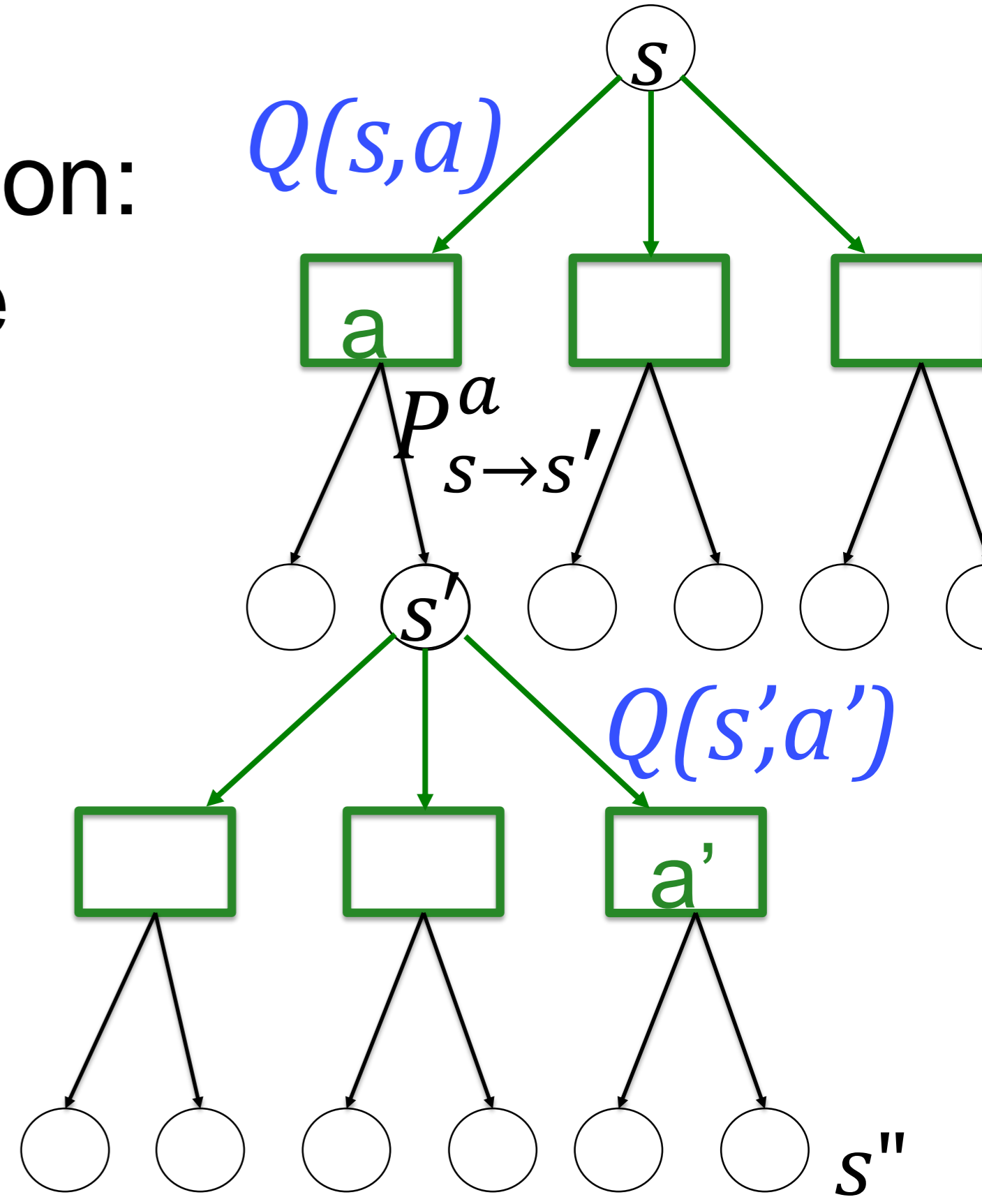
Suppose that we have found a set of Q-values.

We keep them frozen while evaluating expectation:

IFF  $E[\Delta Q(s, a)] = 0$ , then the Q-values solve the Bellman equation.

### Notes:

- Expectation is taken for fixed Q-values and hence for fixed policy.
- Expectation  $E[\Delta Q(s, a)]$  is taken over all possible paths starting in  $(s, a)$ . I call this 'batch-like'.
- The length of the path is given by the needs of the update equation: Here from  $(s, a)$  to  $(r, s', a')$
- Look at state-action-diagram to keep track of terms



Your comments.



# Variant A: SARSA is consistent w. Bellman equation

Blackboard 6A:  
SARSA

We have proven the following:

Suppose that we have found a set of Q-values.

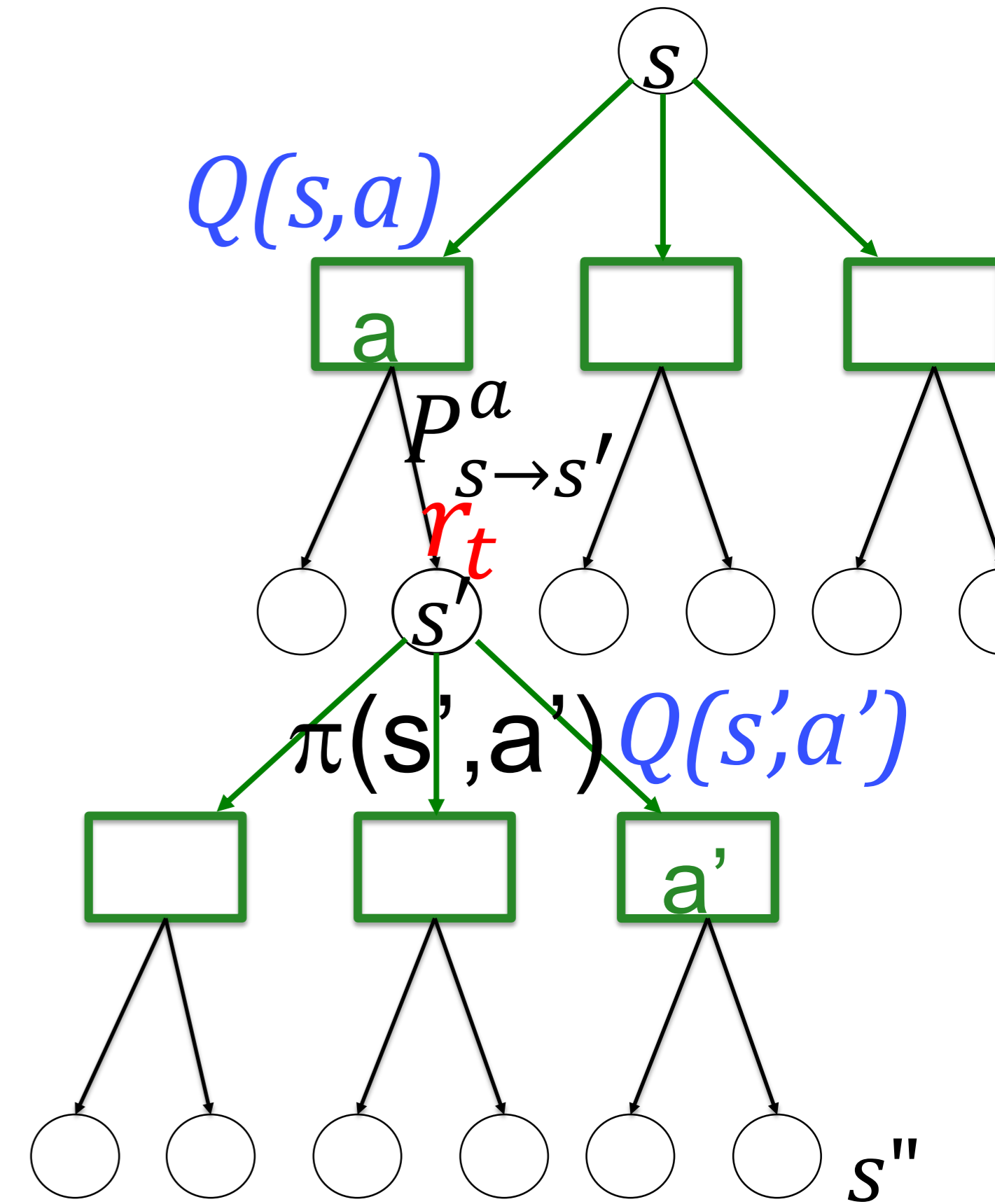
We keep them frozen while evaluating expectation.

IFF  $E[\Delta Q(s, a)] = 0$ , then the Q-values solve the

Bellman equation.

$$E[\Delta Q(s, a)] = \eta E[r_t + \gamma Q(s', a') - Q(s, a)] = 0$$

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



Look at graph to take expectations:

- if algo is on a branch  $(s, a)$ , all remaining expectations are “given  $s$  and  $a$ ”

Previous slide.

This is version A of the theorem. In the proof, we exploit the condition:

$$E[\Delta Q(s, a) | s, a] = 0$$

In order to take the expectations, we look at graph:

- if in the evaluation we are in state  $s'$ , all remaining expectations are “given  $s'$ ”
- if we are on a branch  $(s, a)$ , all remaining exp. are “given  $s$  and  $a$ ”.

We exploit that all Q-values and the policy are fixed while we evaluate the expectation. Hence  $E[Q(s, a)] = Q(s, a)$ .

Note that the proof works in both direction. If the Q-values are those of the Bellman equation, the expected SARSA update step vanishes. And if the expected SARSA update step is zero, then the Q-values correspond to the Bellman equation. Both direction work under the assumption of a fixed policy.

The stronger theorem (with fluctuations, version B) is sketched in the appendix (and video).

Additional Notes: This weaker theorem (variant A that corresponds to the one on the previous slide) takes expectations for FIXED Q-values. We can interpret these expectations as the following ‘batch’ computation

We assume a fixed policy (i.e., under the assumption of a fixed set of Q-values) and a ‘batch version’ of SARSA. Batch-SARSA means that in order to evaluate  $E[\Delta Q(s, a) | s, a]$  we use a large number of starts from the same value (s,a) each time running one step up to (s',a') [note that this gives different (s',a')]. Once the number of starts is large enough to get a full sample of the statistics we update Q(s,a). If the updates with the batch-SARSA do not lead to a change of Q values (for all state-action pairs), then this means that batch-SARSA has converged to the Bellman equation for this fixed policy. (That was the theorem in the main text).

Batch-SARSA is a computational implementation of the way many statistical convergence proofs work: you assume that you average over a full statistical sample of all possibilities given your current state or the current state-action pair. Expectation signs in the update step imply updating over a ‘full batch of data’. In this approach Q-values no longer fluctuate, and hence do not need expectation signs; the policy no longer fluctuates and also does not need expectations signs.

Your comments.

# Variant B: Bellman equation and SARSA: theorem for small $\eta$

Setting: 'temporal averaging'

The SARSA algo has been applied for a very long time, using updates

$$\Delta Q(s, a) = \eta [r_t + \gamma Q(s', a') - Q(s, a)]$$

IF (i) averaged over many update steps

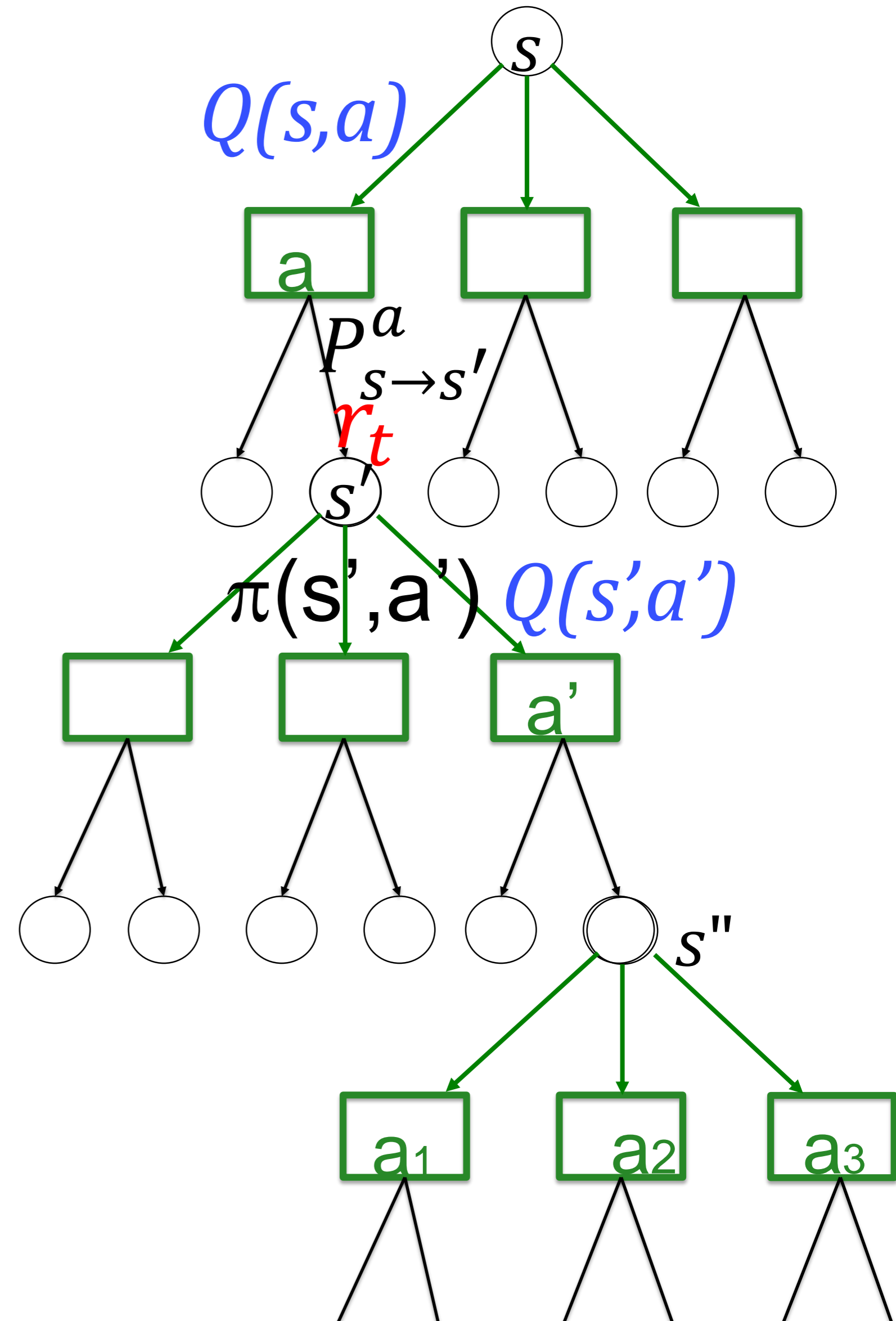
$$\langle \Delta Q(s, a) | s, a \rangle = 0$$

(ii) learning rate  $\eta$  is small ( $\eta \rightarrow 0$ )

THEN: fluctuations of Q are negligible and the set of expected Q-values solves the Bellman eq.

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

with the current policy  $\pi(s', a')$



Previous slide. There are two different versions of the theorem. This version B is proven in the annex, in a fashion similar to the case of the 1-step horizon.

In class (earlier slides) we have shown for the multi-step horizon a weaker statement:

**Expectations over SARSA updates are consistent with the Bellman equation** if the expected update vanishes: **In version A, we assume that Q-values are fixed (frozen) when we take the expectation.**

Note that taking the expectation in version A is different from averaging over update steps in version B. In version A, taking the expectation means that we average over all possible outcomes in the **current** situation, with momentarily fixed Q-values and **fixed policy** (i.e., the one induced by the set of Q-values at time  $t$ ). This distinction is important, because **for a fixed policy averaging is relatively easy.**

However, when averaging over time steps as in variant B, the Q-values and policy are different in each time step, and the proof therefore requires a limit  $\eta \rightarrow 0$ , so that changes can be neglected.

Hence there are therefore two versions of the theorem and two proof-sketches:

Blackboard 6A. On the earlier slides, we assume Q-values are fixed and do not fluctuate.

Blackboard 6B. In the Annex, we assume that Q-values may fluctuate slightly round their stable values. This approach gives additional insights into the situation of the online SARSA, once it has converged in expectation.

# Teaching monitoring – monitoring of understanding

[ ] today, up to here, at least 60% of material was new to me.

[ ] up to here, I have the feeling that I have been able to follow (at least) 80% of the lecture.

# Summary: Reinforcement Learning and SARSA

Now:  
Exercise session

## Learning outcome and conclusions:

- **Reinforcement Learning is learning by rewards**
  - world is full of rewards (but not full of labels)
- **Agents and actions**
  - agent learns by interacting with the environment
  - state  $s$ , action  $a$ , reward  $r$
- **Exploration vs Exploitation**
  - optimal actions are easy if we know reward probabilities
  - since we don't know the probabilities we need to explore
- **Bellman equation**
  - self-consistency condition for Q-values
- **SARSA algorithm: state-action-reward-state-action**
  - update while exploring environment with current policy



# Annex: Variant B - SARSA and Bellman equation (proof for small $\eta$ )

Setting: 'Temporal Averaging'

The SARSA algo with stochastic policy  $\pi$  has been applied for a long time with updates

$$\Delta \hat{Q}(s, a) = \eta [r_t + \gamma \hat{Q}(s', a') - \hat{Q}(s, a)]$$

IF (i) learning rate  $\eta$  is small; AND IF

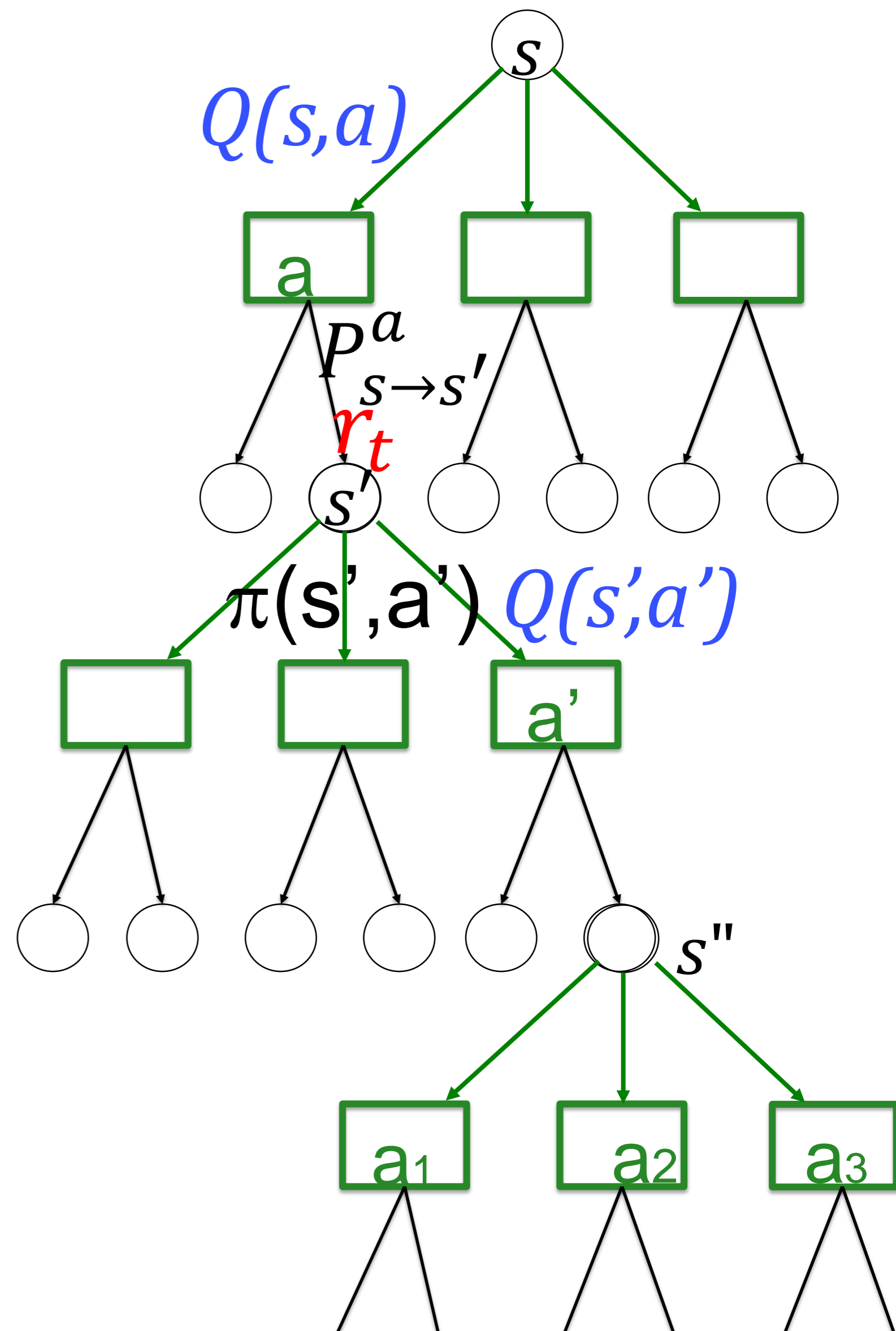
(ii) for all Q-values

$$\langle \Delta Q(s, a) | s, a \rangle = 0$$

THEN the **expectation** values of the set of  $\hat{Q}$ -values solves the Bellman eq.

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

with the current policy  $\pi(s', a')$



Notes: A few points should be stressed:

1. This is not a convergence theorem. We just show consistency as follows:  
if SARSA has converged then it has converged to a solution of the Bellman equation.
2. In fact, for any finite  $\eta$  the SARSA Q-values ( $\hat{Q}$ ) fluctuate a little bit. It is only the EXPECTATION value of the  $\hat{Q}$  which converges.
3. We should keep in mind that SARSA is an on-policy online algorithm for arbitrary state-transition graphs. Hence the value  $\hat{Q}$  at  $(s,a)$  and  $(s',a')$  will both fluctuate!
4. The policy depends on these  $\hat{Q}$ -values and hence fluctuates as well.  
To keep fluctuations of the policy small, we need small  $\eta$ .  
We imagine that all  $\hat{Q}$  values fluctuate around their expectation value with small standard deviation. As a result,  $\pi$  also fluctuates around a 'standard' policy.
5. The fluctuations of the policy can be smaller than that of the Q-values: for example in epsilon-greedy, you first order actions by the value of  $Q(s,a)$ , and then only the rank of  $Q(s,a)$  matters, not their exact values. In the proof we assume that the fluctuations of the policy become negligible ( $\rightarrow$  shift policy outside expectation).
6. We show that the Q-values in the sense of Bellman are the expectation values of the  $\hat{Q}$  in the sense of SARSA.
7. Expectations are over many trials of the ONLINE SARSA.

The statement and proof is different to slide 127 and to the book of Sutton and Barto.

Blackboard6:  
SARSA

# SARSA is consistent with Bellman ( $\eta \rightarrow 0$ )

$$\frac{1}{\eta} \Delta \hat{Q}(s, a) = \underbrace{\eta}_{\eta} [r_t + \gamma \hat{Q}(s', a') - \hat{Q}(s, a)]$$

$$0 = \frac{1}{\eta} E(\Delta \hat{Q}(s, a)) = E[r_t] + \gamma E[\hat{Q}(s', a')] - E[\hat{Q}(s, a)]$$

$$E[\hat{Q}(s, a)] \text{ given } (s, a) = \sum_{s'} P_{s \rightarrow s'}^a R_{s \rightarrow s'}^a + \gamma \sum_{s'} P_{s \rightarrow s'}^a \sum_{a'} E[\pi(s', a') \cdot \hat{Q}(s', a')]$$

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

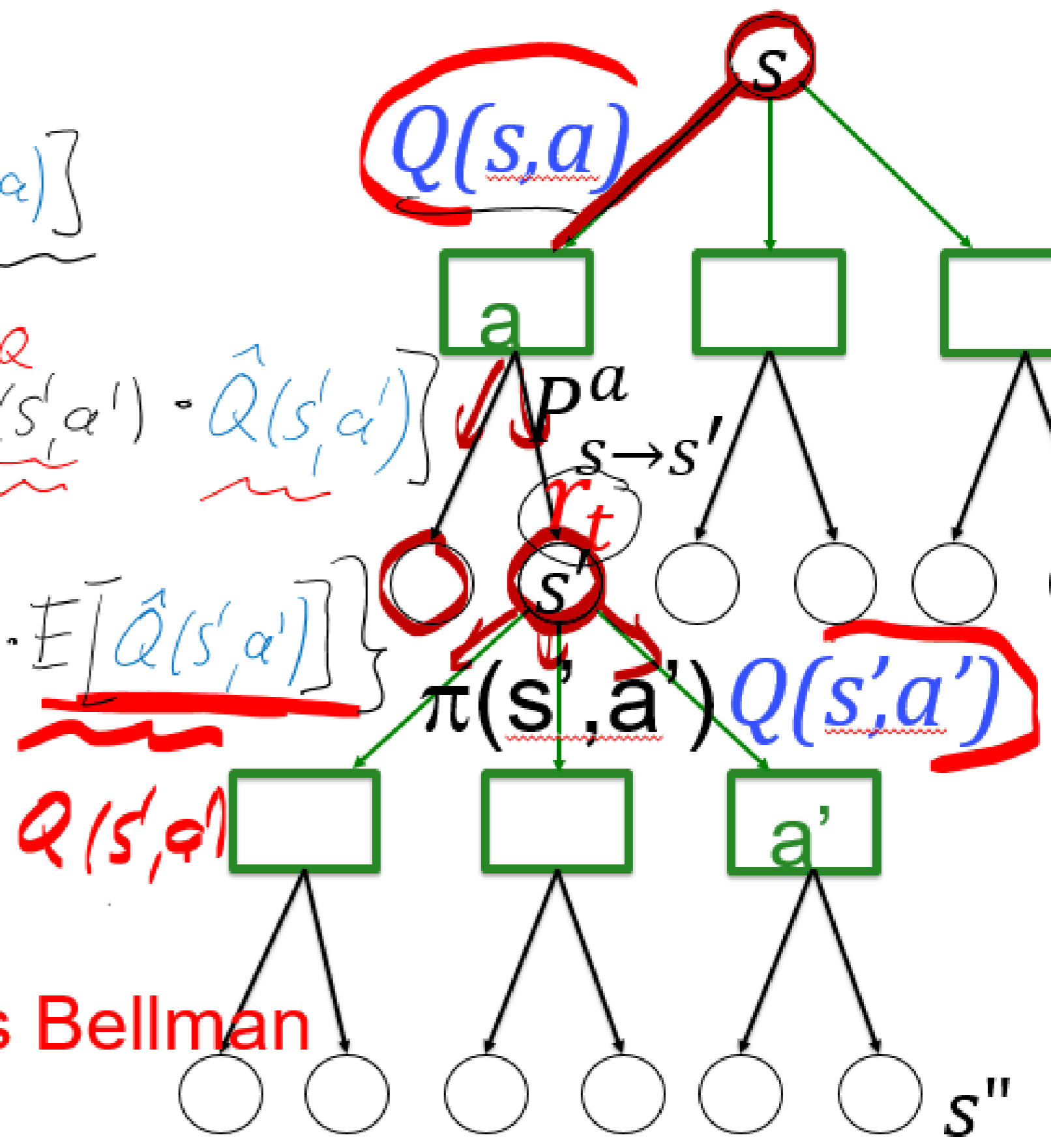
$Q(s, a)$

Bellman

Online:  $Q(s, a)$  jitters, but its mean (expectation) solves Bellman

Look at graph to take expectations:

- if algo is in state  $s$ , all remaining expectations are "given  $s$ "
- if algo is on a branch  $(s, a)$ , all remaining exp. are "given  $s$  and  $a$ "



Notes: This is a proof sketch of the consistency of online SARSA (**Variant B**). We allow all Q-values to fluctuate around their expectation (visualized as ‘temporal averaging’), but we still have to keep fluctuations of the policy negligibly small.

If we allow for small fluctuations of the policy, then we have to realize that these fluctuations are correlated with the fluctuations of Q-values. Thus the evaluation of the product  $E(\pi Q)$  is not trivial. Moreover, correlations can lead to a shift of the value and make the result inconsistent with the Bellman equation.

This variant B therefore only works in the limit of vanishing learning rate.

The other theorem (**Variant A** that corresponds to the one on the slides in the main part of this lecture) takes expectations for FIXED Q-values. We can interpret these expectations as corresponding to a ‘batch’ computation.

Batch-SARSA is a computational implementation of the way many statistical convergence proofs work: you assume that you average over a full statistical sample of all possibilities given your current state or the current state-action pair. Expectation signs in the update step imply updating over a ‘full batch of data’. In this approach Q-values no longer fluctuate, and hence do not need expectation signs; the policy no longer fluctuates and also does not need expectations signs.