# Lecture 09: Model free learning in RL

**Radoslav Neychev**

# Outline

1. MDP properties recap
2. Value-function and Q-function
3. Reward discounting
4. Q-learning, temporal difference
5. Approximate Q-learning

Based on: https://github.com/yandexdataschool/Practical_RL/ weeks 2, 3 and 4
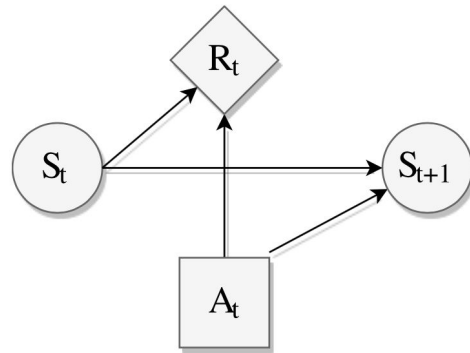
# References

These slides are deeply based on [Practical RL course](#) weeks 2, 3 and 4
Special thanks to YSDA team for making them publicly available.

# Given dynamics, how to find an optimal policy?

**Definition of Markov Decision Process**

MDP is a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$, where

1. $\mathcal{S}$ – set of states of the world
2. $\mathcal{A}$ – set of actions
3. $\mathcal{P} : \mathcal{S} \times \mathcal{A} \mapsto \triangle(\mathcal{S})$ – state-transition function, giving us $p(s_{t+1} \mid s_t, a_t)$
4. $\mathcal{R} : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ – reward function, giving us $\mathbb{E}_R\left[R(s_t, a_t) \mid s_t, a_t\right]$.



**Markov property**

$$p(r_t, s_{t+1} \mid s_0, a_0, r_0, ..., s_t, a_t) = p(r_t, s_{t+1} \mid s_t, a_t)$$

(next state, expected reward) depend on (previous state, action)

# Goal: solve an MDP by finding an optimal policy

1.  What is the objective?
    a.  Reward: discounting and design
    b.  Expected objective: state- and action-value function

# Explaining goals to agent through reward

Reward hypothesis (R.Sutton)

Goals and purposes can be thought of as the maximization of the expected value of the cumulative sum of a received scalar signal

**Cumulative reward** is called a return:

end of an episode

$$G_t \triangleq R_t + R_{t+1} + R_{t+2} + ... + R_T$$

immediate reward

E.g.: reward in **chess** – value of taken opponent's piece

# E.g.: data center non-stop cooling system

- **S**tates – temperature measurements
- **A**ctions – different fans speed
- **R = 0**  for exceeding temperature thresholds
- **R = +1** for each second system is cool

## What could go wrong with such a design?

# E.g.: data center non-stop cooling system

- **S**tates – temperature measurements
- **A**ctions – different fans speed
- **R = 0** for exceeding temperature thresholds
- **R = +1** for each second system is cool

What could go wrong with such a design?

Infinite return for non optimal behaviour!

$$G_t = 1 + 1 + 0 + 1 + 1 + 0 + .... = \sum_{t=1}^{\infty} R_t = \infty$$

# E.g.: cleaning robot

- **S**tates – dust sensors, air
- **A**ctions – cleaning / rest / conditioning on or off
- **R = 100**  for long tedious floor cleaning task done
- **R = 1**  for turning air conditioning on-off
- Episode ends each day

## What could go wrong with such a design?

OpenAI blog post about faulty rewards: https://openai.com/blog/faulty-reward-functions/
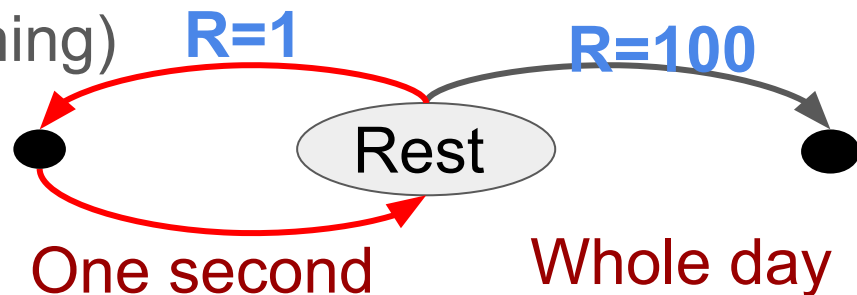
# E.g.: cleaning robot

- **S**tates – dust sensors, air
- **A**ctions – cleaning / rest / conditioning on or off
- **R = 100**  for long tedious floor cleaning task done
- **R = 1**  for turning air conditioning on-off
- Episode ends each day

## What could go wrong with such a design?

Reward(air)  <  Reward(cleaning)
Time(air)  <<  Time(cleaning)

**R=1**          **R=100**

Rest

Positive feedback loop!

One second          Whole day

# Explaining goals to agent through reward

Reward hypothesis (R.Sutton)

Goals and purposes can be thought of as the maximization of the expected value of the cumulative sum of a received scalar signal

Get rid of infinite sum by <span style="color:red">discounting</span>    $0 \le \gamma < 1$

$$G_t \triangleq R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + ... = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

<span style="color:red">discount factor</span>

The same cake compared to today's one worth

$\gamma$

- $\gamma^2$ times less tomorrow

- times less the day after

tomorrow



$\gamma$ <span style="color:red">will eat it day by day</span>

# Reward discounting

## Maximal return for **R = +1**

$$G_0 = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

| $\gamma$ | 0.9 | 0.95 | 0.99 |
|----------|-----|------|------|
| $\frac{1}{1-\gamma}$ | 10 | 20 | 100 |

R=+1 discounted $n$ times
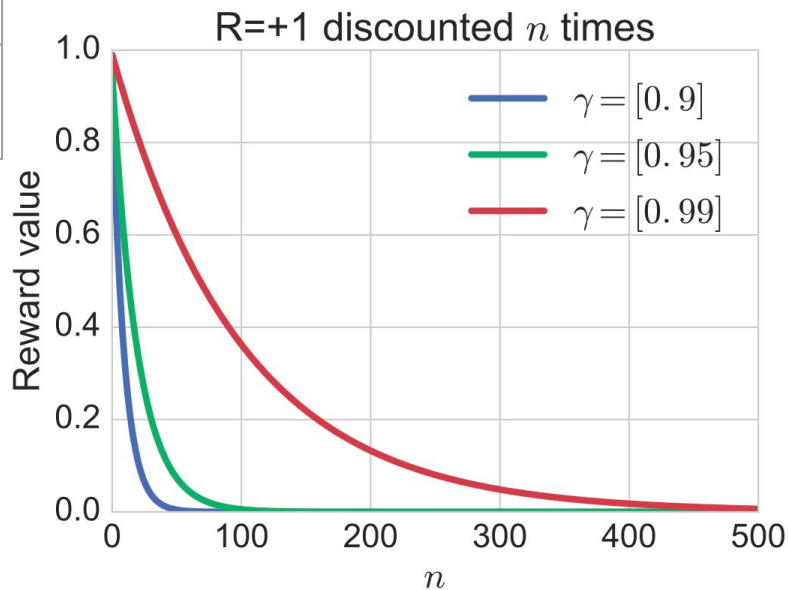
# Discounting makes sums finite

Maximal return for **R = +1**

$$G_0 = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

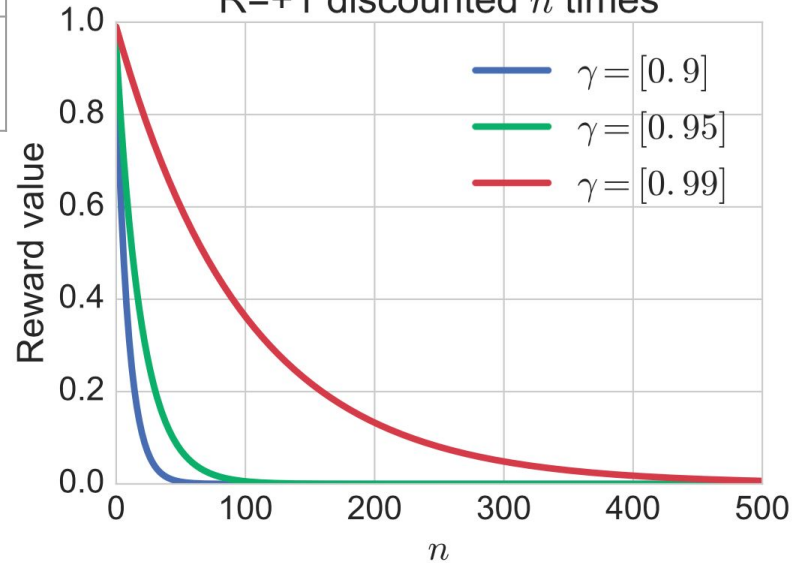| $\gamma$ | 0.9 | 0.95 | 0.99 |
|---|---|---|---|
| $\frac{1}{1-\gamma}$ | 10 | 20 | 100 |

Any discounting changes optimisation task and its solution!



R=+1 discounted $n$ times

# **State-** and **Action-**value functions

v(s) is <u>expected</u> return conditional on <u>state</u>:

$$v_\pi(s) \triangleq \mathbb{E}_\pi \left[ G_t \,|\, S_t = s \right]$$

$$= \mathbb{E}_\pi \left[ R_t + \gamma G_{t+1} \,|\, S_t = s \right]$$

$$= \sum_a \pi(a \,|\, s) \sum_{r,s'} p(r, s' \,|\, s, a) \left[ r + \gamma \boxed{\mathbb{E}_\pi \left[ G_{t+1} | S_{t+1} = s' \right]} \right]$$

$$= \sum_a \pi(a \,|\, s) \sum_{r,s'} p(r, s' \,|\, s, a) \left[ r + \gamma \underline{v_\pi(s')} \right]$$

<span style="color:red">By definition</span>

<span style="color:green">Intuition</span>: value of following policy $\pi$ from state s

# **Action-value** function q(s, a)

Is <u>expected</u> <span style="color:darkred">return</span> conditional on <u>state</u> **and action**:

Intuition: value of following policy $\pi$ <u>after</u> committing action **a** in state **s**

$$
\begin{aligned}
q_\pi(s, a) &= \mathbb{E}_\pi \left[ G_t \mid S_t = s, A_t = a \right] \\
&= \mathbb{E}_\pi \left[ R_t + \gamma G_{t+1} \mid S_t = s, A_t = a \right] \\
&= \sum_{r, s'} p(r, s' \mid s, a) \left[ r + \gamma \mathbb{E}_\pi \left[ G_{t+1} \mid S_{t+1} = s' \right] \right] \\
&= \sum_{r, s'} p(r, s' \mid s, a) \left[ r + \gamma v_\pi(s') \right]
\end{aligned}
$$

# Relations between v(s) and q(s,a)

We already know how to write q(s,a) in terms of v(s)

$$q_\pi(s, a) = \sum_{r, s'} p(r, s' \,|\, s, a) \left[ r + \gamma v_\pi(s') \right]$$

What about v(s) in terms of q(s,a)?

$$v_\pi(s) = \sum_a \pi(a \,|\, s) \sum_{r, s'} p(r, s' \,|\, s, a) \left[ r + \gamma v_\pi(s') \right]$$

$$= \sum_a \pi(a \,|\, s) q_\pi(s, a)$$

So, we could now write q(s, a) in terms of q(s,a)!

$$q_\pi(s, a) = \sum_{r, s'} p(r, s' \,|\, s, a) \left[ r + \gamma \sum_{a'} \pi(a' \,|\, s') q_\pi(s', a') \right]$$

# Bellman expectation equation for v(s)

Recursive definition of v(s) is an important concept in RL

$$v_\pi(s) = \sum_a \pi(a \mid s) \sum_{r,s'} p(r, s' \mid s, a) \left[ r + \gamma v_\pi(s') \right]$$

$$= \mathbb{E}_\pi \left[ R_t + \gamma v_\pi(S_{t+1}) \mid S_t = s \right]$$
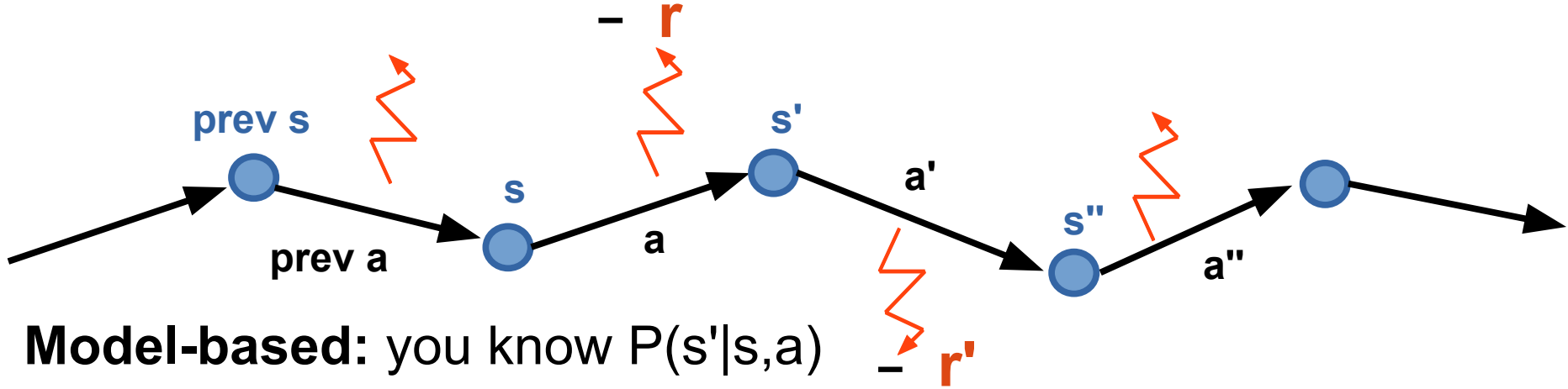
# Bellman expectation equation for **q(s,a)**

$$q_\pi(s, a) = \sum_{r, s'} p(r, s' \mid s, a) \left[ r + \gamma v_\pi(s') \right]$$

$$= \sum_{r, s'} p(r, s' \mid s, a) \left[ r + \gamma \sum_{a'} \pi(a' \mid s') q_\pi(s', a') \right]$$

To sum up

- **Vπ(s)** – expected G from state **s** if you follow **π**

- **V*(s)** – expected G from state s if you follow **π*** – optimal

- **Qπ(s,a)** – expected G from state **s**

  – if you start by taking action **a**

  – and follow **π** from next state on

- **Q*(s,a)** – same as Qπ(s,a) where **π** = **π*** – optimal policy

$$Q^*(s,a) = \underset{s',r}{E}\, r(s,a) + \gamma \cdot V^*(s')$$
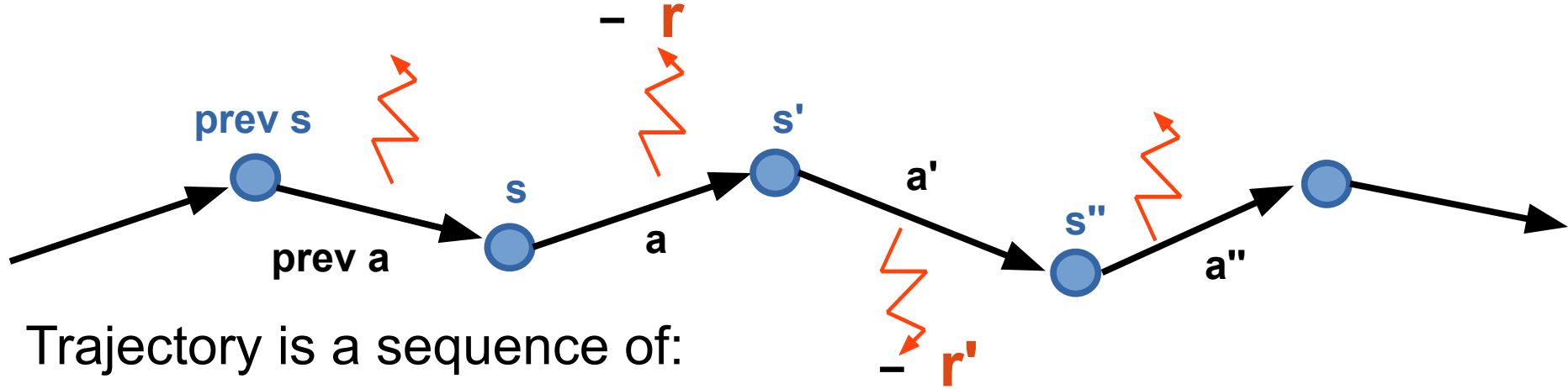
$$V^*(s) = max\, Q^*(s,a)$$

# Learning from trajectories



**Model-based:** you know P(s'|s,a)
- can apply dynamic programming
- can plan ahead

**Model-free:** you can sample trajectories
- can try stuff out
- insurance not included

# Learning from trajectories

Trajectory is a sequence of:
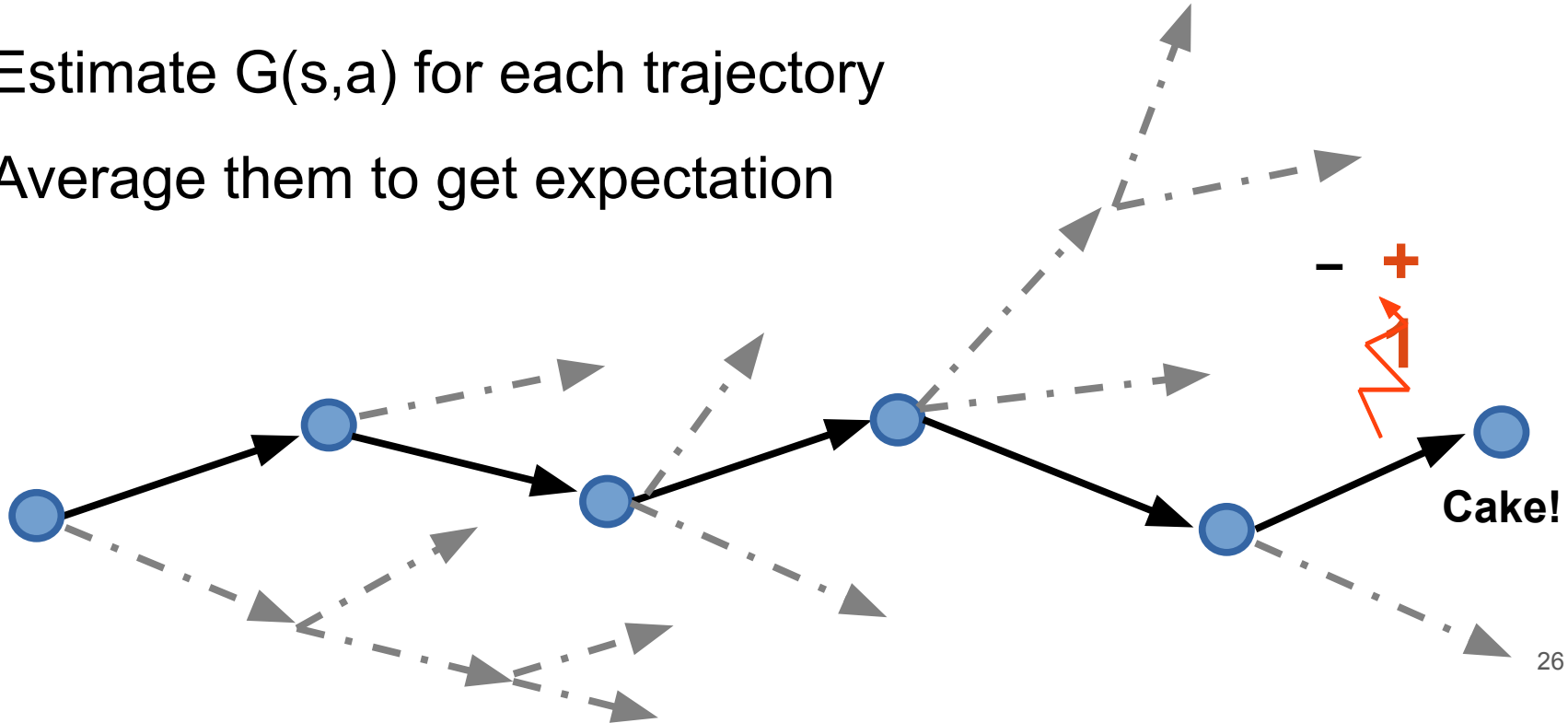- states (s)
- actions (a)
- rewards (r)

We can only sample trajectories

**Q:** What to learn?
V(s) or Q(s,a)

V(s) is useless
without P(s'|s,a)

- Get all trajectories containing particular (s,a)

- Estimate G(s,a) for each trajectory
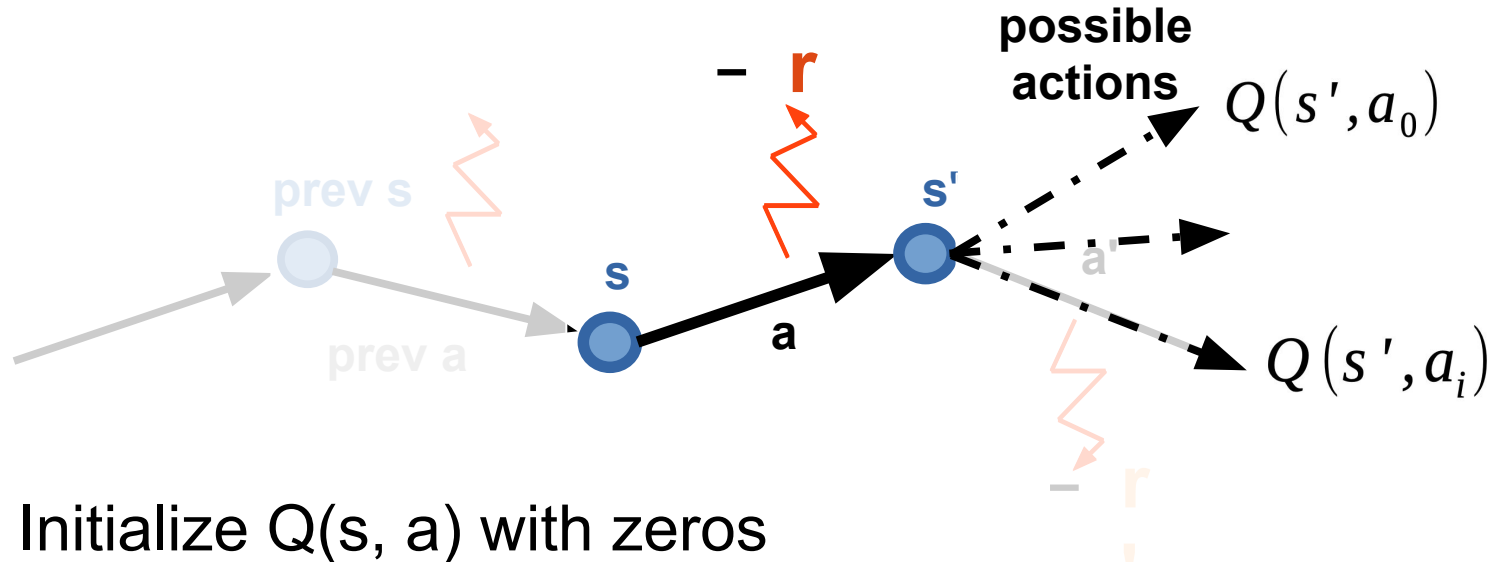
- Average them to get expectation

− **+**

**Cake!**

- Q(s, a) can be improved iteratively!

$$Q(s_t, a_t) \leftarrow \underset{r_t, s_{t+1}}{E} \; r_t + \gamma \cdot max_{a'} Q(s_{t+1}, a')$$

**That's Q\*(s,a)**

**That's value for π\***
aka optimal policy

**That's something
we don't have**

**What do we do?**

- Q(s, a) can be improved iteratively!

$$Q(s_t, a_t) \leftarrow \underset{r_t, s_{t+1}}{E} \, r_t + \gamma \cdot max_{a'} \, Q(s_{t+1}, a')$$

$$\underset{r_t, s_{t+1}}{E} \, r_t + \gamma \cdot max_{a'} Q(s_{t+1}, a') \approx \frac{1}{N} \sum_i r_i + \gamma \cdot max_{a'} Q(s_i^{next}, a')$$

$$Q(s_t, a_t) \leftarrow \alpha \cdot (r_t + \gamma \cdot max_{a'} Q(s_{t+1}, a')) + (1 - \alpha) Q(s_t, a_t)$$

Initialize Q(s, a) with zeros
- Sample <s, a, r, s'> from the environment
- Compute new Q(s, a) eslimation:

$$\hat{Q}(s,a)=r(s,a)+\gamma\,\underset{a_i}{max}\,Q(s',a_i)$$

- Update Q(s, a):

$$Q(s,a)\leftarrow\alpha\cdot\hat{Q}(s,a)+(1-\alpha)Q(s,a)$$

29

# Monte-carlo

Averages Q over
sampled paths



**+1** **Cake!**

# Temporal Difference

Uses recurrent
formula for Q



**possible actions**

_− r_

**s** **a** **s'**

$$Q^*(s,a) = E_{s',r} \, r(s,a) + \gamma \cdot V^*(s')$$

$$Q(s_t, a_t) \leftarrow \alpha \cdot (r_t + \gamma \cdot max_{a'} Q(s_{t+1}, a')) + (1-\alpha) Q(s_t, a_t)$$

$$\pi(s): argmax_a Q(s,a)$$

Strategies:

- ε-greedy
  - With probability ε take random action; otherwise take optimal action.

- Softmax
  Pick action proportional to softmax of shifted normalized Q-values.

$$\pi(a|s) = softmax\left(\frac{Q(s,a)}{\tau}\right)$$

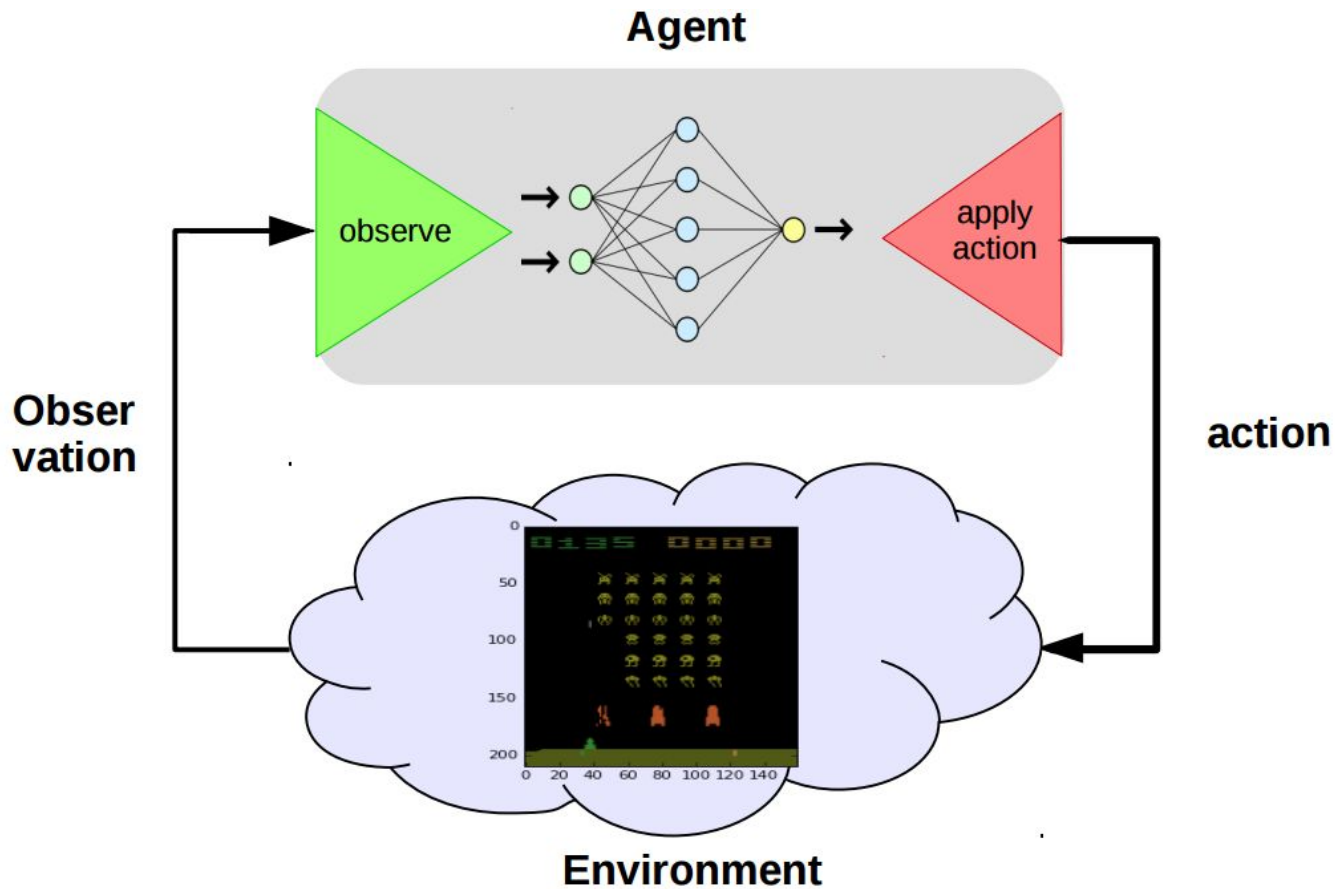For now states and actions are discrete. What if states are **not**?

Minimize loss given <**s, a, r, s'**>

$$L = [Q(s_t, a_t) - Q^{true}(s_t, a_t)]^2$$

$$L \approx [Q(s_t, a_t) - (r_t + \gamma \cdot max_{a'} Q(s_{t+1}, a'))]^2$$

**Agent**

observe

apply action

**Obser vation**

**action**

**Environment**

$$G_t = \sum_{t'=t}^{T} \gamma^{(t'-t)} r_{t'}$$

$$Q^{\pi}(s, a) = E_{\pi}[G_t | s_t = s, a_t = a]$$

$$V^{\pi}(s) = E_{\pi}[G_t | s_t = s]$$

**Recurrent relations**

$$Q^{\pi}(s, a) = E_{s_{t+1}}[r_t + \gamma V^{\pi}(s_{t+1})]$$

$$Q^{\pi}(s, a) = E_{s_{t+1}, a_{t+1} \sim \pi}[r_t + \gamma Q^{\pi}(s_{t+1}, a_{t+1})]$$

For all $\pi, s, a$:      $Q^{\pi^*}(s, a) \geq Q^{\pi}(s, a)$

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

**Bellman optimality equation**

$$Q^*(s_t, a) = E_{s_{t+1}}[r_t + \max_{a'} Q^*(s_{t+1}, a')]$$

**Training step**

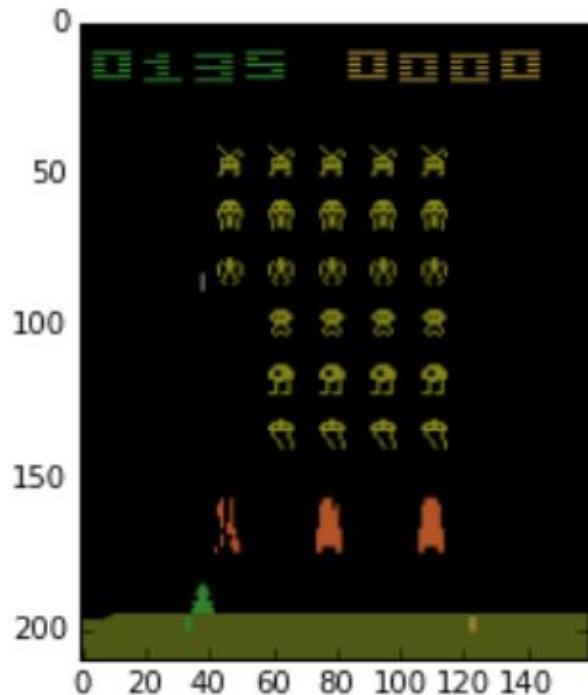$$Q(s_t, a_t) \longleftarrow Q(s_t, a_t) + \alpha(r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t))$$

**Q-learning as MSE minimization**

$$L = (r_t + \gamma \max_{a'} \boxed{Q(s_{t+1}, a')} - Q(s_t, a_t))^2$$

Const

$$\nabla L = 2 \cdot (r_t + \gamma \max_{a'} \boxed{Q(s_{t+1}, a')} - Q(s_t, a_t))$$

**What's wrong here?**

How many states are there?

approximately
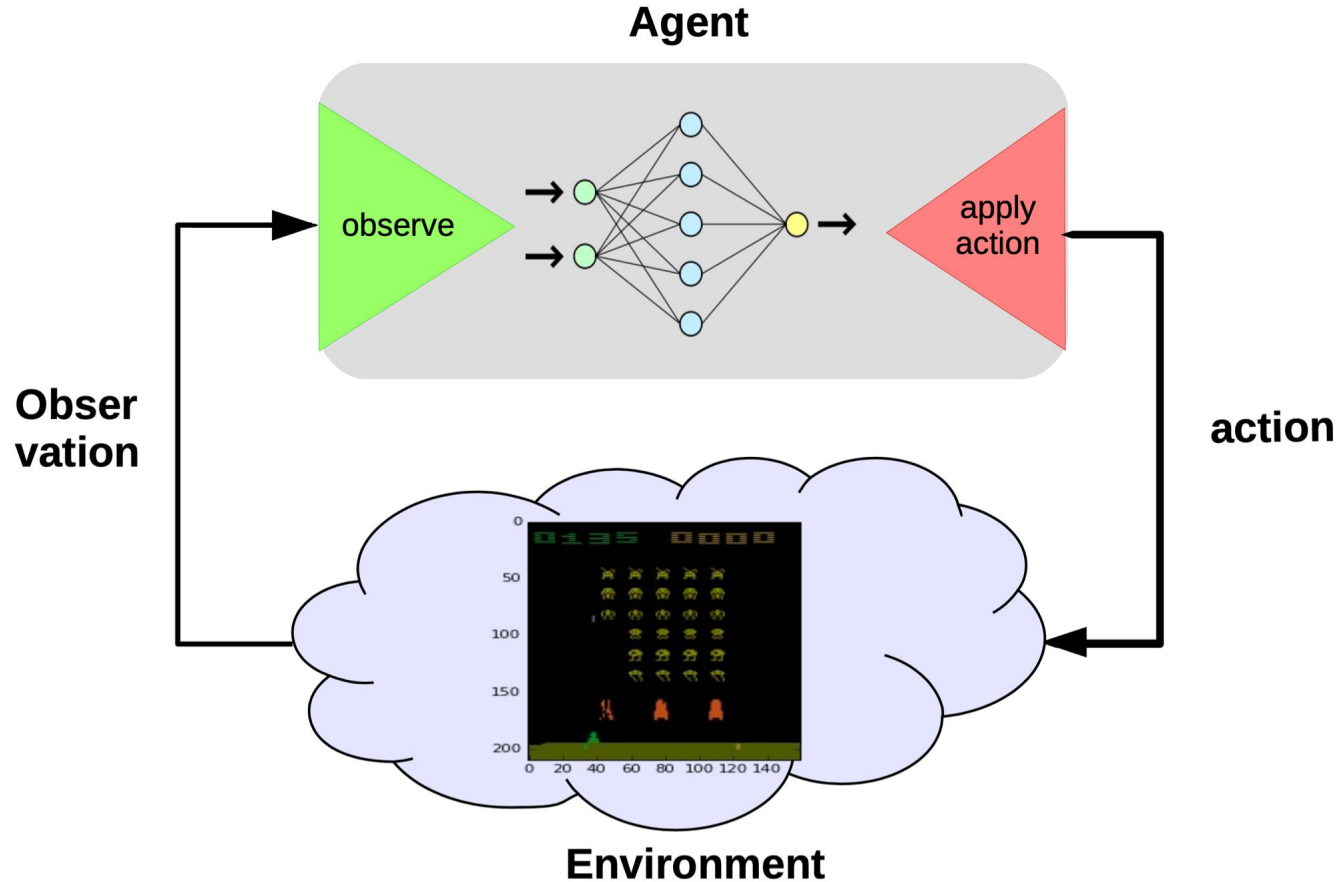
$$|S| = 2^{210 \cdot 160 \cdot 8 \cdot 3}$$

For now states and actions are discrete. What if states are **not**?

Minimize loss given <**s, a, r, s'**>

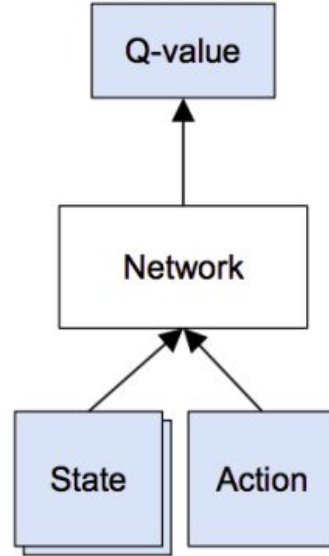$$L = \left[ Q(s_t, a_t) - Q^{true}(s_t, a_t) \right]^2$$

$$L \approx \left[ Q(s_t, a_t) - (r_t + \gamma \cdot max_{a'} Q(s_{t+1}, a')) \right]^2$$
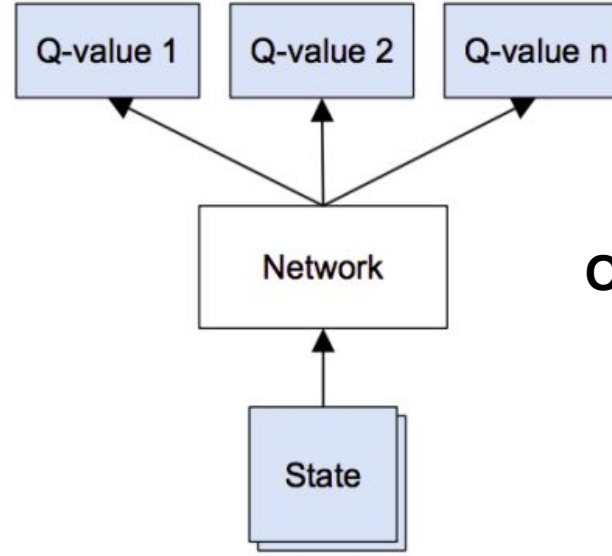
# Approximate Q-learning

# Possible architectures
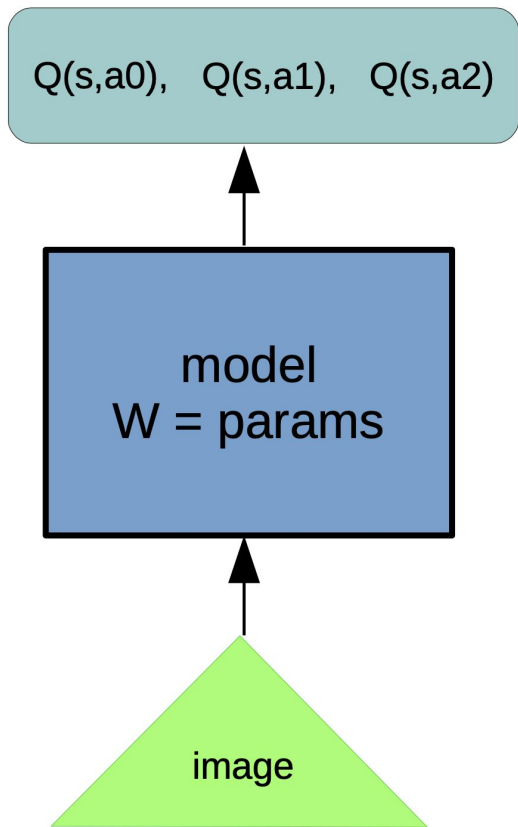


**Continuous control or large number of actions**

**One pass for all actions**

Q-value

Network

State     Action

Given **(s,a)**
Predict Q(s,a)

Q-value 1    Q-value 2    Q-value n

Network

State

Given **s** predict all q-values
Q(s,a0), Q(s,a1), Q(s,a2)
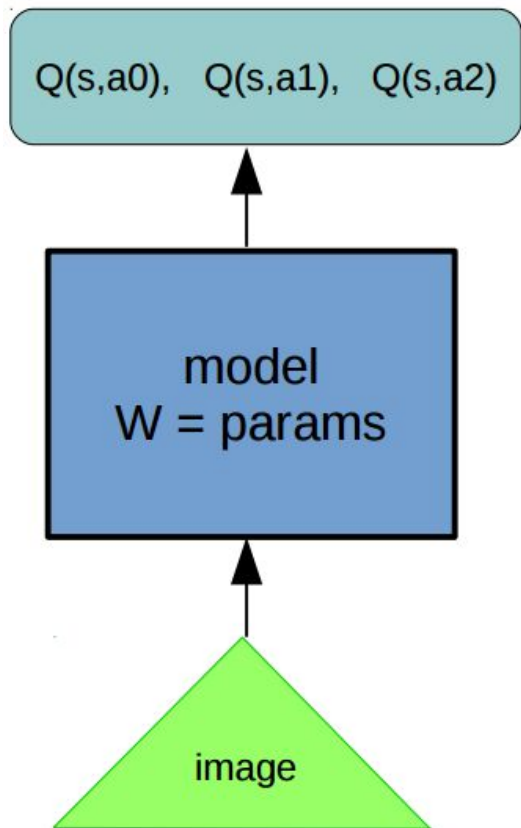
Q(s,a0),  Q(s,a1),  Q(s,a2)

model
W = params

image

$$\hat{Q}(s_t, a_t) = r + \gamma \cdot max_{a'} \hat{Q}(s_{t+1}, a')$$

$$L = \left(Q(s_t, a_t) - \left[r + \gamma \cdot max_{a'} Q(s_{t+1}, a')\right]\right)^2$$

Consider const

$$w_{t+1} = w_t - \alpha \cdot \frac{\delta L}{\delta w}$$

# Approximate Q-learning

Q(s,a0),  Q(s,a1),  Q(s,a2)

model
W = params

image

**Objective:**

$$L=\left(Q\left(s_t,a_t\right)-\hat{Q}\left(s_t,a_t\right)\right)^2$$

**consider const**

**Q-learning:**

$$\hat{Q}\left(s_t,a_t\right)=r+\gamma\cdot max_{a'}\,Q\left(s_{t+1},a'\right)$$

**SARSA:**

$$\hat{Q}\left(s_t,a_t\right)=r+\gamma\cdot Q\left(s_{t+1},a_{t+1}\right)$$

**Expected Value SARSA:**

$$\hat{Q}\left(s_t,a_t\right)=r+\gamma\cdot \underset{a'\sim\pi(a|s)}{E}\,Q\left(s_{t+1},a'\right)$$

What kind of network digests images well?

# Basic deep Q-learning



Qvalues → action → apply action

Qvalues is a dense layer with **no** nonlinearity

Dense

Conv

Conv

Observation

**ε-greedy** rule (tune ε or use probabilistic rule)

**Whatever** you found in your favorite deep learning toolkit

- Q-learning allows to learn some approximation of the reward function and environment model
  - So we can use it to solve the desired problem


- Remember what Q(s, a) and V(s) functions do
- Remember both about exploration and exploitation
  - At least using greedy policy or softmax smoothing