

Lecture 2.0

Reinforcement Learning: Q-learning

Alexey Gruzdev
alexey.s.gruzdev@gmail.com
HSE, Winter 2020

Recap: Dynamic Programming

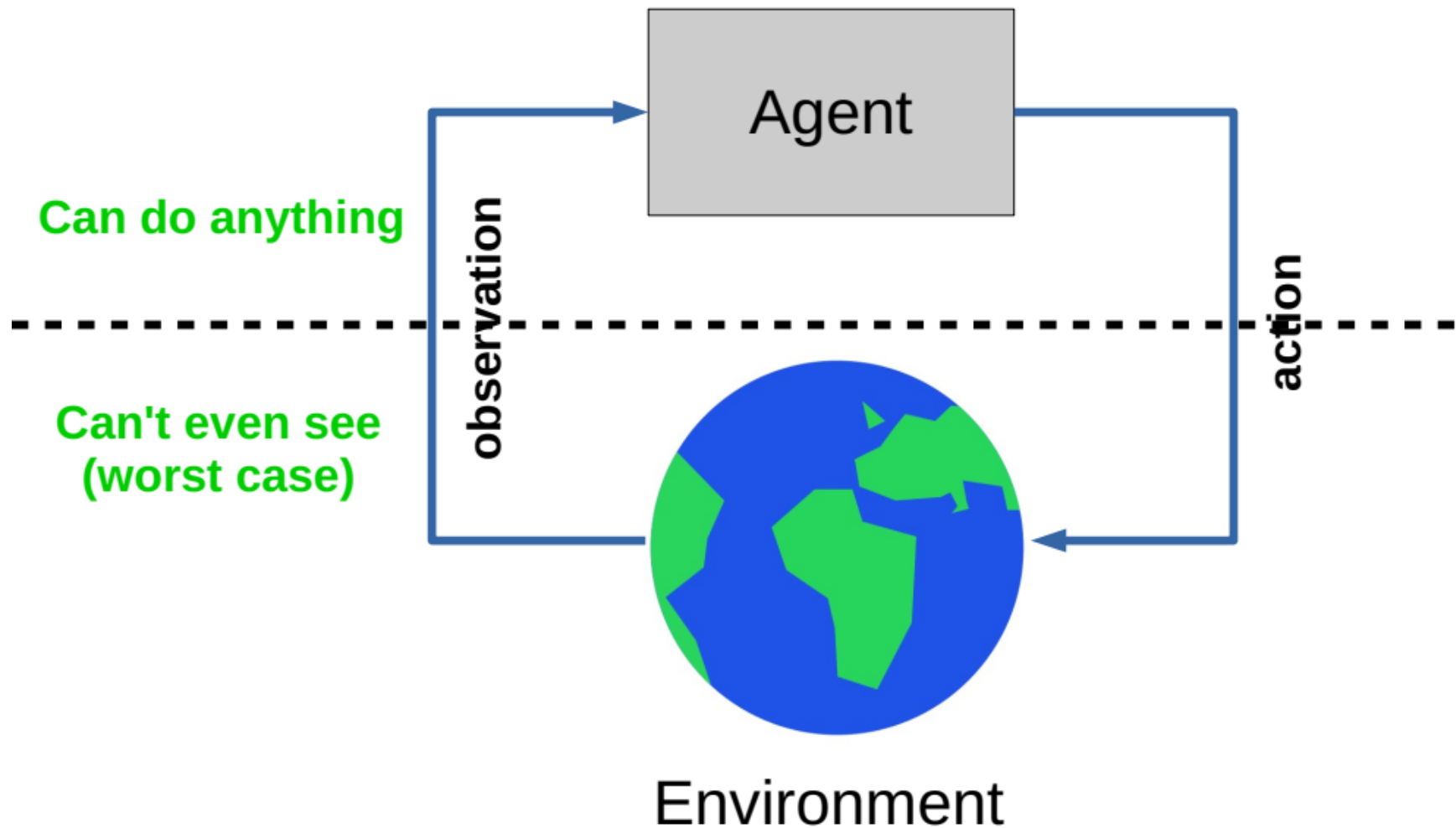
- $v_{\pi}(s), v_*(s)$
- If you know $v_*(s), p(r, s' | s, a) \rightarrow$ know optimal policy
- We can learn $v_*(s)$ with Dynamic Programming:

$$v_*(s) = \max_a \sum_{r,s'} p(r, s' | s, a) [r + \gamma v_*(s')]$$

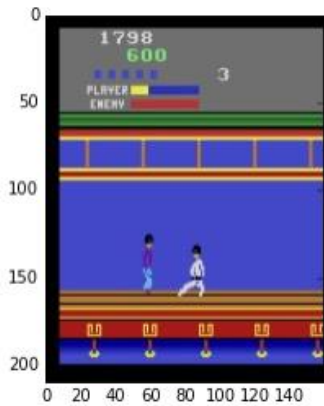
- $q_{\pi}(s, a), q_*(s, a)$

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

Decision making: reality check



Decision making: reality check



Model-Free Setup

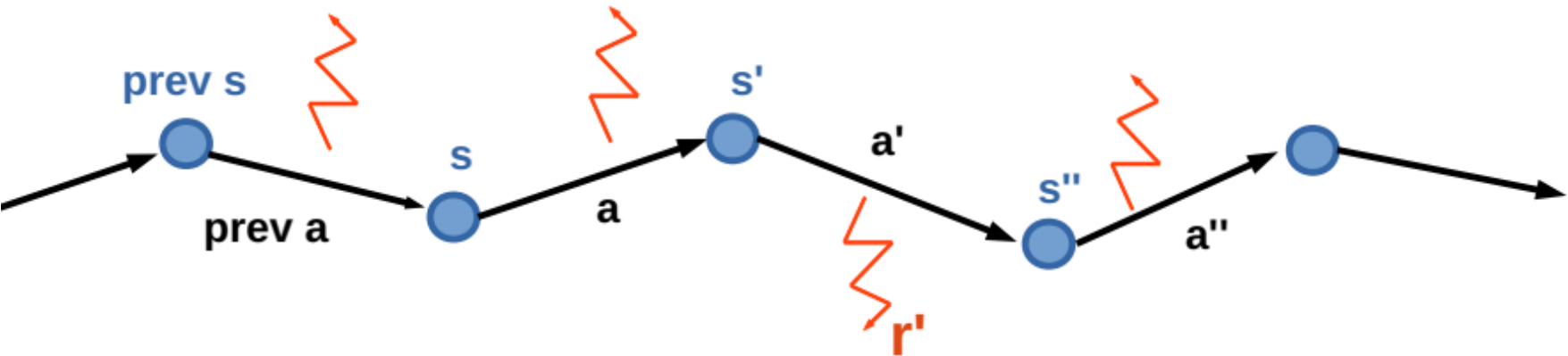
- We don't know internal environment representation, e.g.

$$p(r, s' \mid s, a) - \text{unknown}$$

What should we do?

Learning from trajectories

- $s_1 \rightarrow a_1 \rightarrow r_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_n$ – trajectory



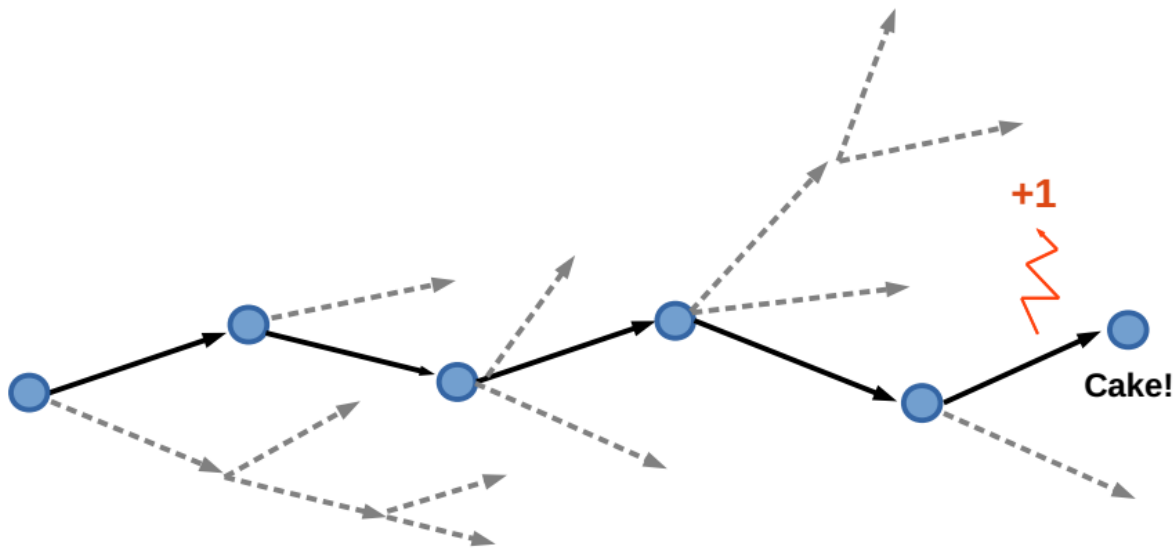
- Model-based setup:
 - you can apply Dynamic Programming
 - you can plan (!)
- Model-free setup:
 - you can experiment with different actions
 - no guaranties (!!!)

Learning from trajectories

- $s_1 \rightarrow a_1 \rightarrow r_1 \rightarrow s_2 \rightarrow \dots \rightarrow s_n$ – trajectory
- We can sample trajectories (a lot of trajectories!)
- What should we learn ?
 - $p(r, s' \mid s, a)$
 - $v_\pi(s)$
 - $q_\pi(s, a)$

Monte-Carlo RL

- Just like N+1 heuristic:
 - Get all trajectories containing particular (s, a)
 - Estimate $G_t(s, a)$ for each trajectory
 - Average them to get *estimation* of expectation



Monte-Carlo RL

- MC methods learn directly from episodes of experience
- MC is *model-free*: no knowledge of MDP transitions / rewards
- MC learns from *complete* episodes: no bootstrapping
- MC uses the simplest possible idea: value = mean return
- Note: can only apply MC to *episodic* MDPs
 - All episodes must terminate

Temporal Difference

- Just like in the 'incremental mean' example we can improve $q_{\pi}(s, a)$ iteratively:

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

- We don't have $p(r, s' | s, a)$ to compute 'fair' expectation, so what should we do?

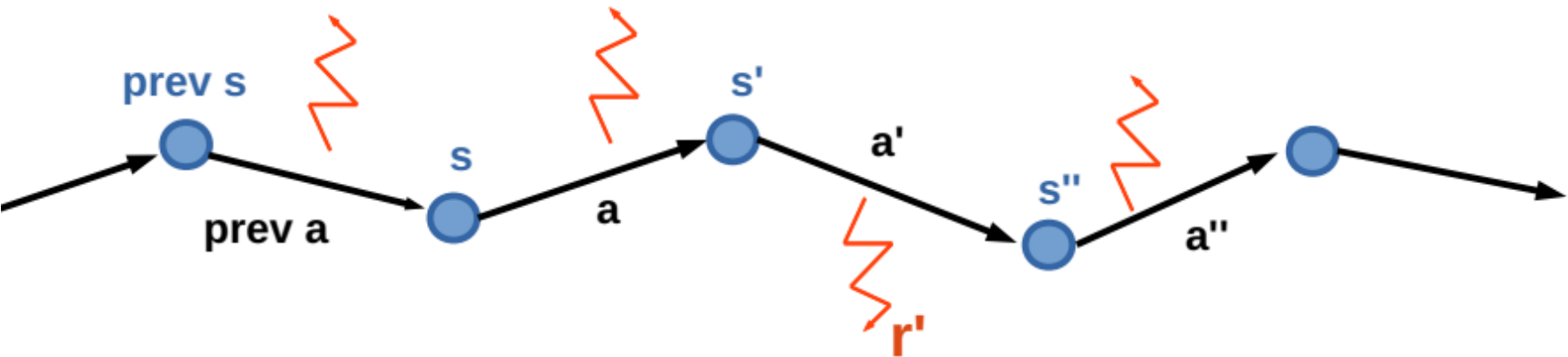
Temporal Difference

$$\begin{aligned} \sum_{r,s'} p(r, s' \mid s, a) [r + \gamma \max_{a'} q_*(s', a')] &\approx \\ \approx \frac{1}{N} \sum_i r_i + \gamma \max_{a'} Q(s'_i, a') \end{aligned}$$

- One more trick: use alpha-smoothing for updating Q-values.

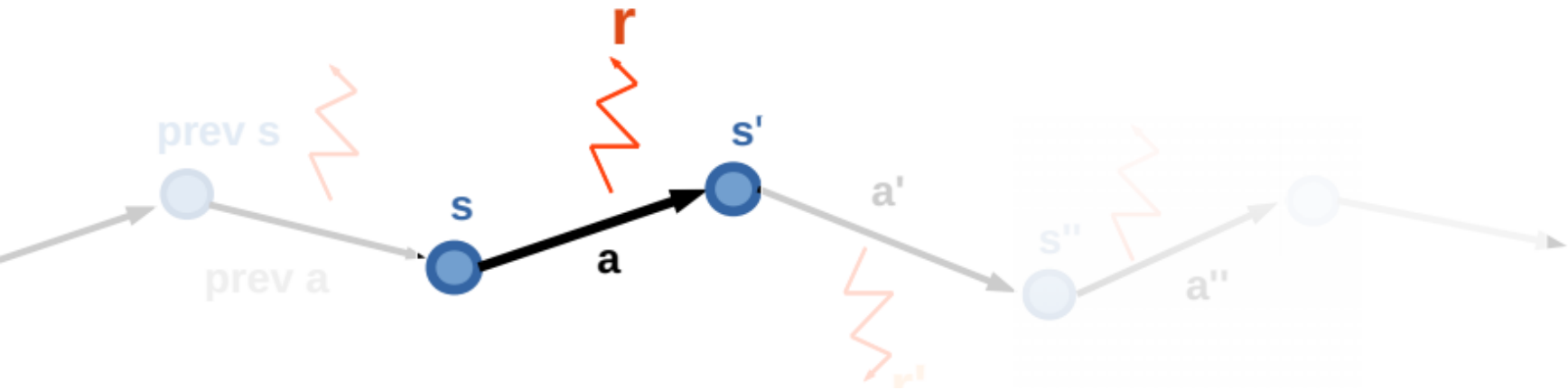
$$Q(s_t, a_t) = \alpha (r_t + \gamma \max_{a'} Q(s_{t+1}, a')) + (1 - \alpha) Q(s_t, a_t)$$

Q-learning



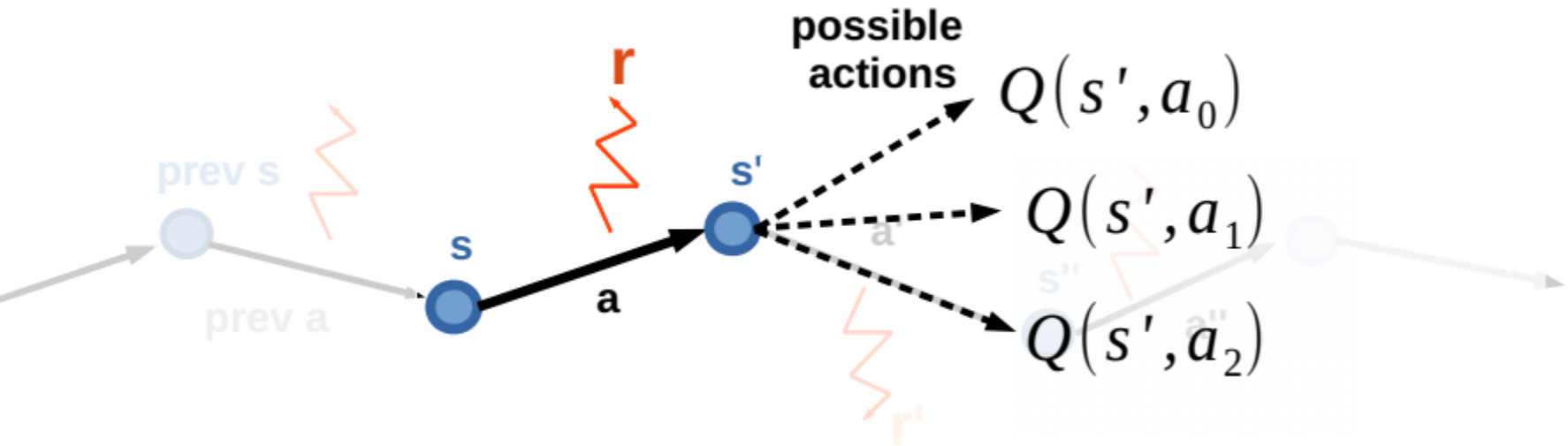
- Works on a sequence of
 - states (s)
 - actions (a)
 - rewards (r)

Q-learning



- Initialize $Q(s, a)$ with zeros
- Cycle:
 - Sample $\langle s, a, r, s' \rangle$ from environment

Q-learning



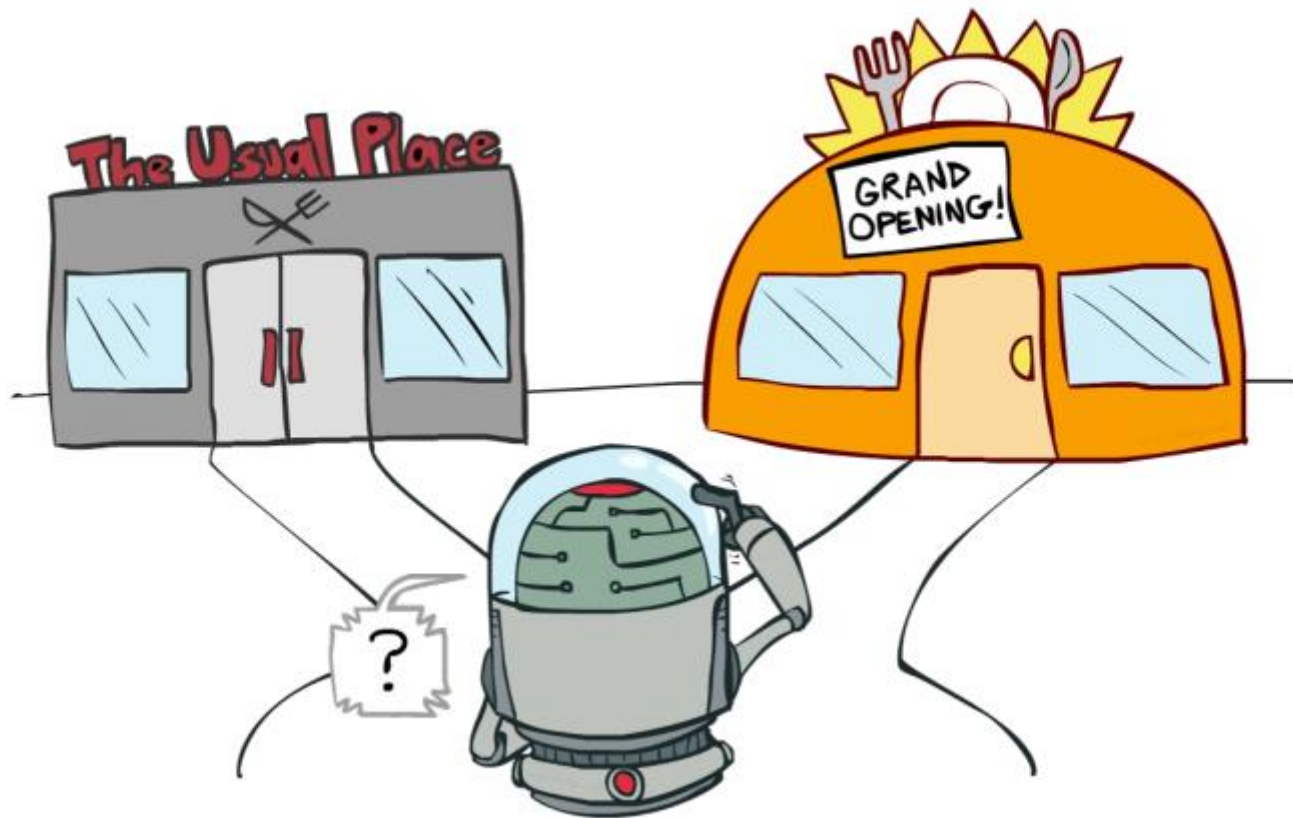
- Initialize $Q(s, a)$ with zeros
- Cycle:
 - Sample $\langle s, a, r, s' \rangle$ from environment
 - Compute $\hat{Q}(s, a) = r(s, a) + \gamma \max_{a_i} Q(s', a_i)$
 - Update: $Q(s_t, a_t) = \alpha \hat{Q}(s, a) + (1 - \alpha) Q(s_t, a_t)$

MC vs TD

- TD can learn *before* knowing the final outcome
 - TD can learn online after every step
 - MC must wait until end of episode before return is known
- TD can learn *without* the final outcome
 - TD can learn from incomplete sequences
 - MC can only learn from complete sequences
 - TD works in continuing (non-terminating) environments
 - MC only works for episodic (terminating) environments

Exploration/Exploitation Revisited

- Balance between using what you learned and trying to find something even better



Exploration/Exploitation Revisited

- Strategies:
 - ϵ -greedy
With probability ϵ take random action,
otherwise take optimal action.
 - ϵ - dithering
Adding random noise to Q-values with ϵ probability

Exploration/Exploitation over time

- If you want to converge to optimal policy you need to gradually reduce exploration.
- Example:
Initialize ϵ -greedy $\epsilon = 0.5$, then gradually reduce it
 - If $\epsilon \rightarrow 0$, it's **greedy in the limit**
 - Be careful with non-stationary environments

Temporal Difference Learning

- TD methods learn directly from episodes of experience
- TD is *model-free*: no knowledge of MDP transitions / rewards
- TD learns from *incomplete* episodes, by *bootstrapping*
- TD updates a guess towards a guess

Reinforcement Learning in the Wild

- Reinforcement learning can be used to solve *large* problems, e.g.
 - Backgammon: **1020 states**
 - Computer Go: **10170 states**
 - Helicopter: **continuous state space**
- How can we scale up the model-free methods for *prediction* and *control* ?

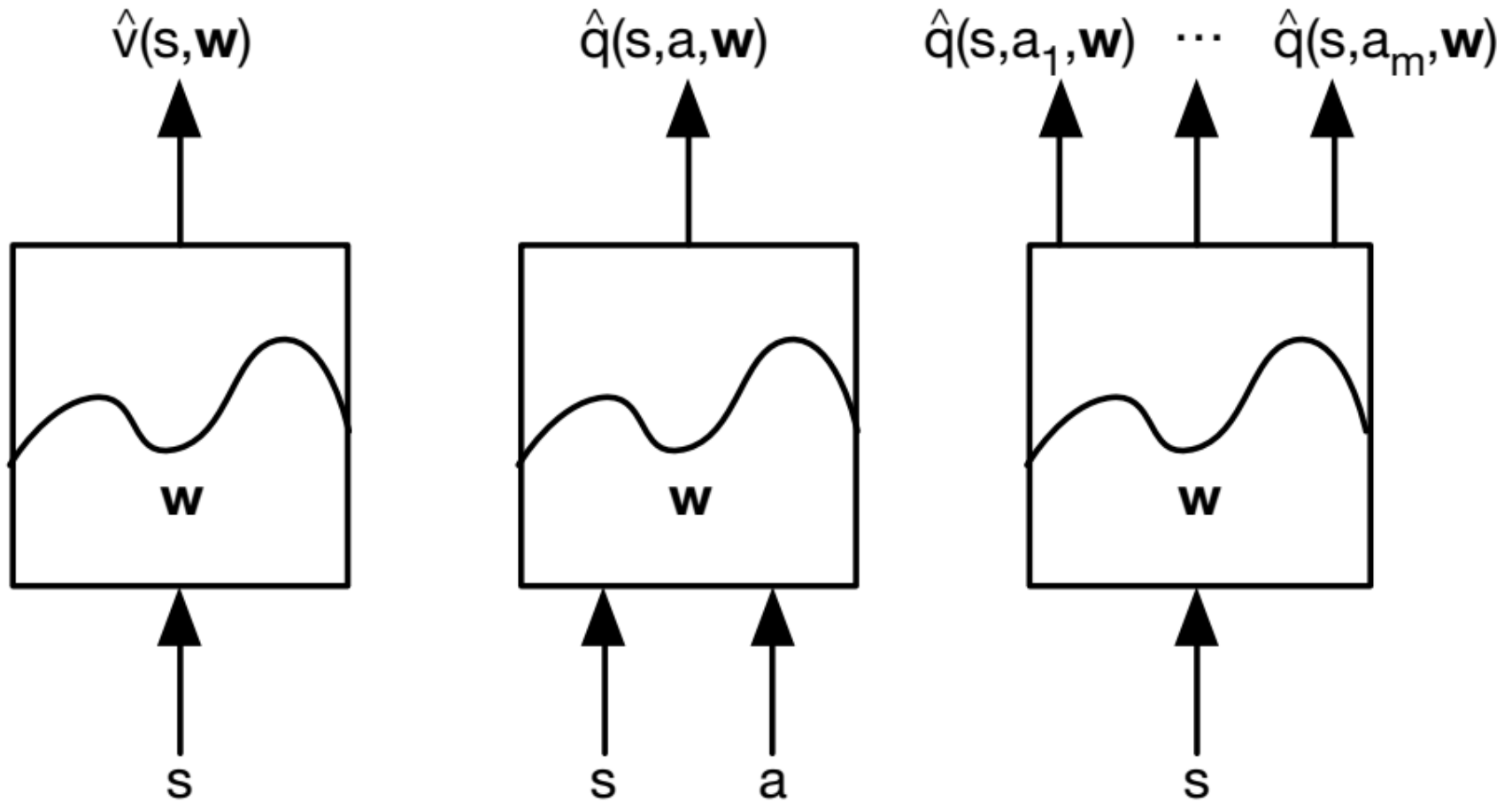
Curse of dimensionality in RL

Problem:

- State space is usually large, sometimes continuous.
- How about action space ?
- However, states do have a structure, similar states have similar action outcomes

What should we do?

Types of Value Function Approximation



Which class of function to choose?

- There are many function approximators, e.g.
 - *Linear combinations of features*
 - *Neural network*
 - Decision tree
 - Nearest neighbor
 - Fourier / wavelet bases
 - ...

Gradient Descent

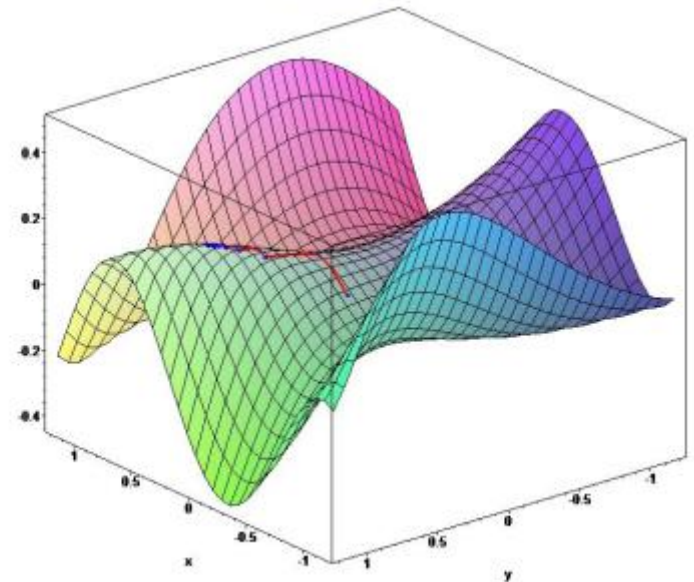
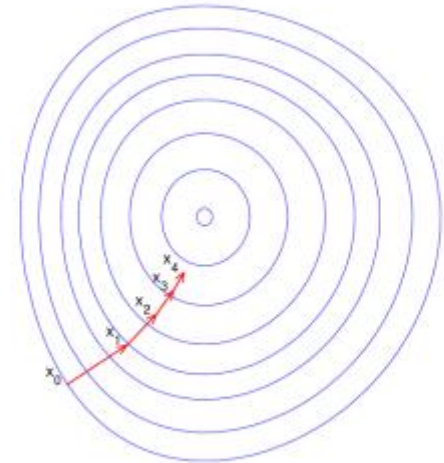
- Let $J(\mathbf{w})$ be a differentiable function of parameter vector \mathbf{w}
- Define the *gradient* of $J(\mathbf{w})$ to be

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{pmatrix} \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}_n} \end{pmatrix}$$

- To find a local minimum of $J(\mathbf{w})$:
 - Adjust \mathbf{w} in direction of -ve gradient

$$\Delta \mathbf{w} = -\frac{1}{2} \alpha \nabla_{\mathbf{w}} J(\mathbf{w})$$

where α is a step-size parameter



SGD for Value Function approximation

- Goal: find parameter vector \mathbf{w} minimizing mean-squared error between approximate value $v'(s, \mathbf{w})$ and true value $v_\pi(s)$

$$J(\mathbf{w}) = \mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, \mathbf{w}))^2]$$

- Gradient descent finds a local minimum

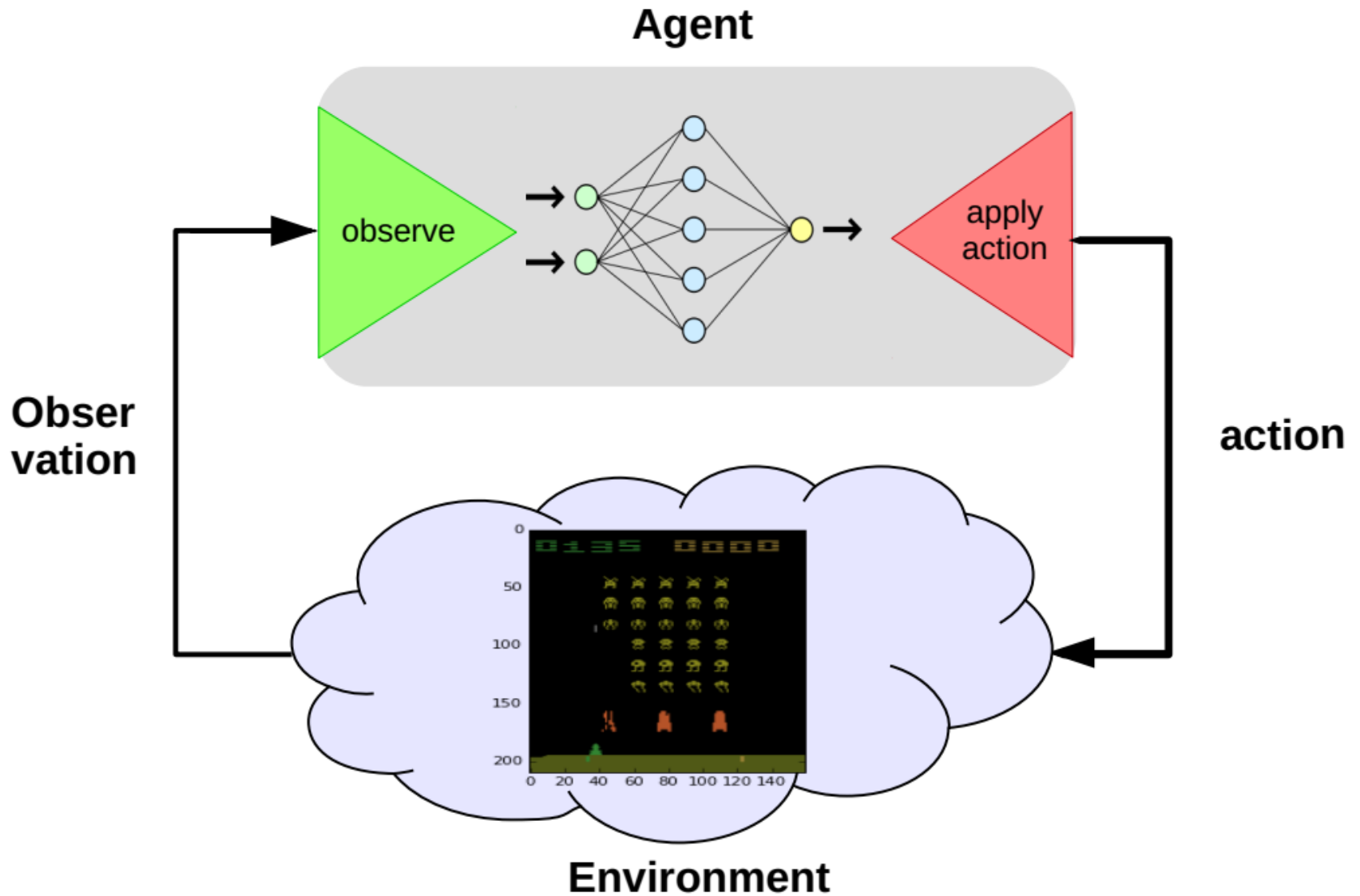
$$\begin{aligned}\Delta \mathbf{w} &= -\frac{1}{2}\alpha \nabla_{\mathbf{w}} J(\mathbf{w}) \\ &= \alpha \mathbb{E}_\pi [(v_\pi(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})]\end{aligned}$$

- Stochastic gradient descent *samples* the gradient

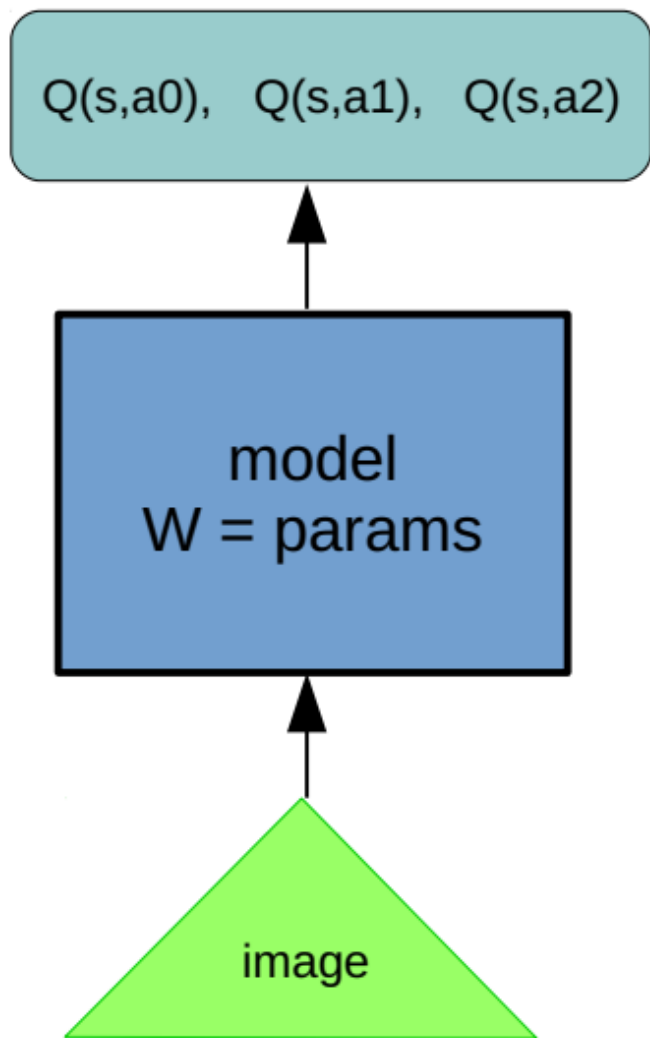
$$\Delta \mathbf{w} = \alpha (v_\pi(S) - \hat{v}(S, \mathbf{w})) \nabla_{\mathbf{w}} \hat{v}(S, \mathbf{w})$$

- Expected update is equal to full gradient update

Atari again



Approximate Q-learning

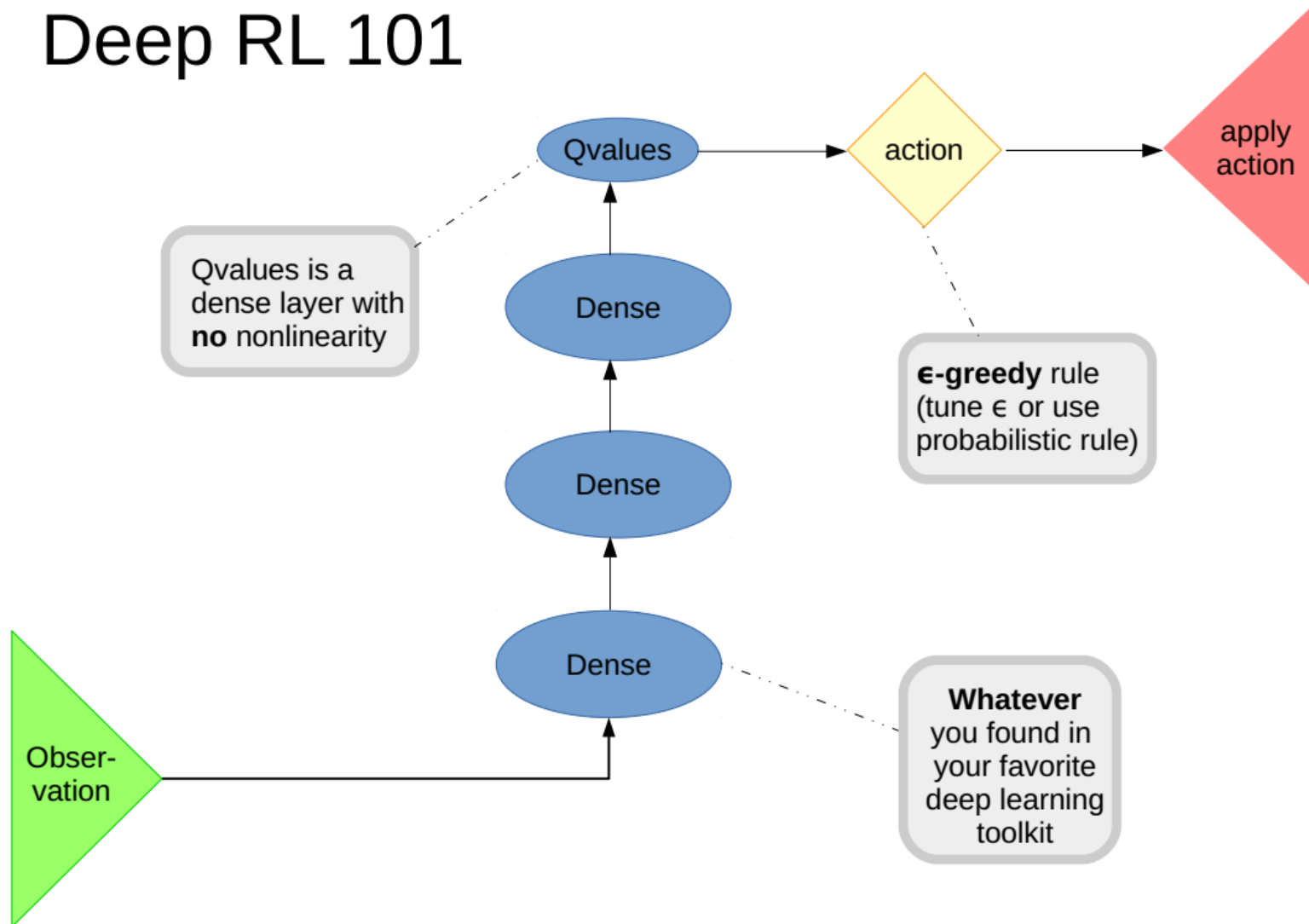


- Initialize W .
- Cycle:
 - Sample $\langle s, a, r, s' \rangle$ from environment
 - Compute $\hat{Q}(s, a) = r(s, a) + \gamma \max_{a_i} Q(s', a_i)$
 - Objective:
$$L = [Q(s_t, at) - \hat{Q}(s_t, at)]^2$$
 - SGD Update:

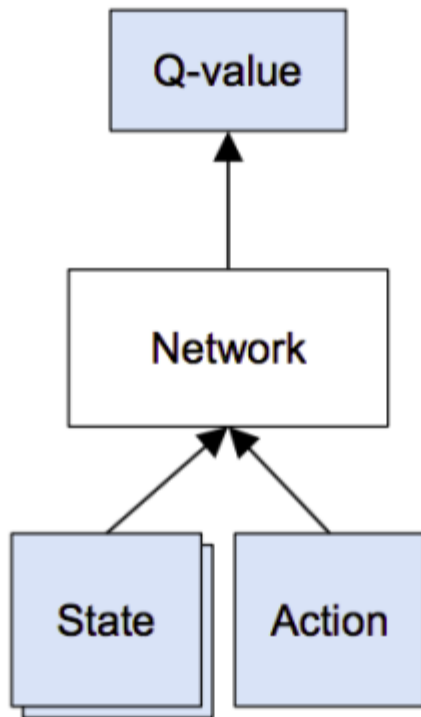
$$W_{t+1} = W_t - \alpha \frac{\partial L}{\partial w_t}$$

RL Mechanics

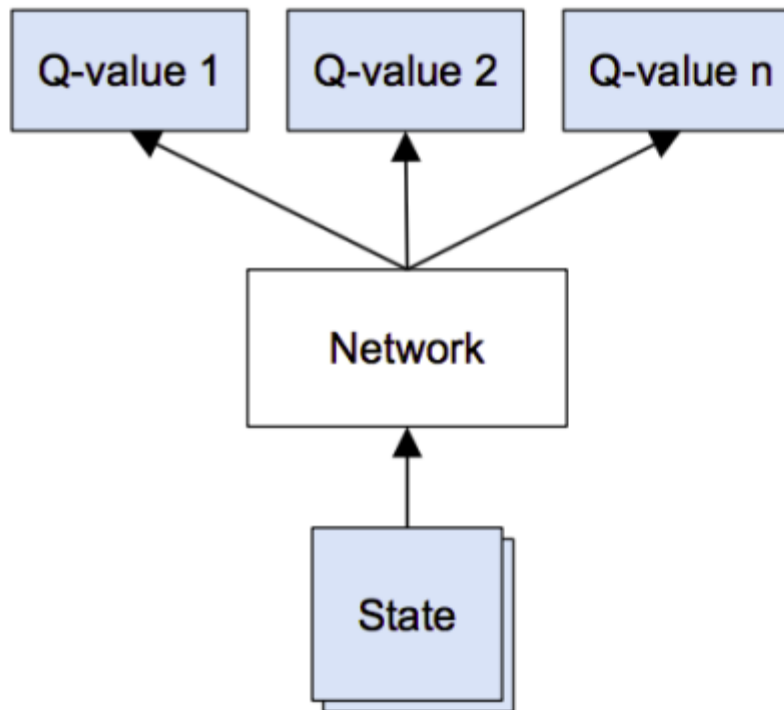
Deep RL 101



Architectures

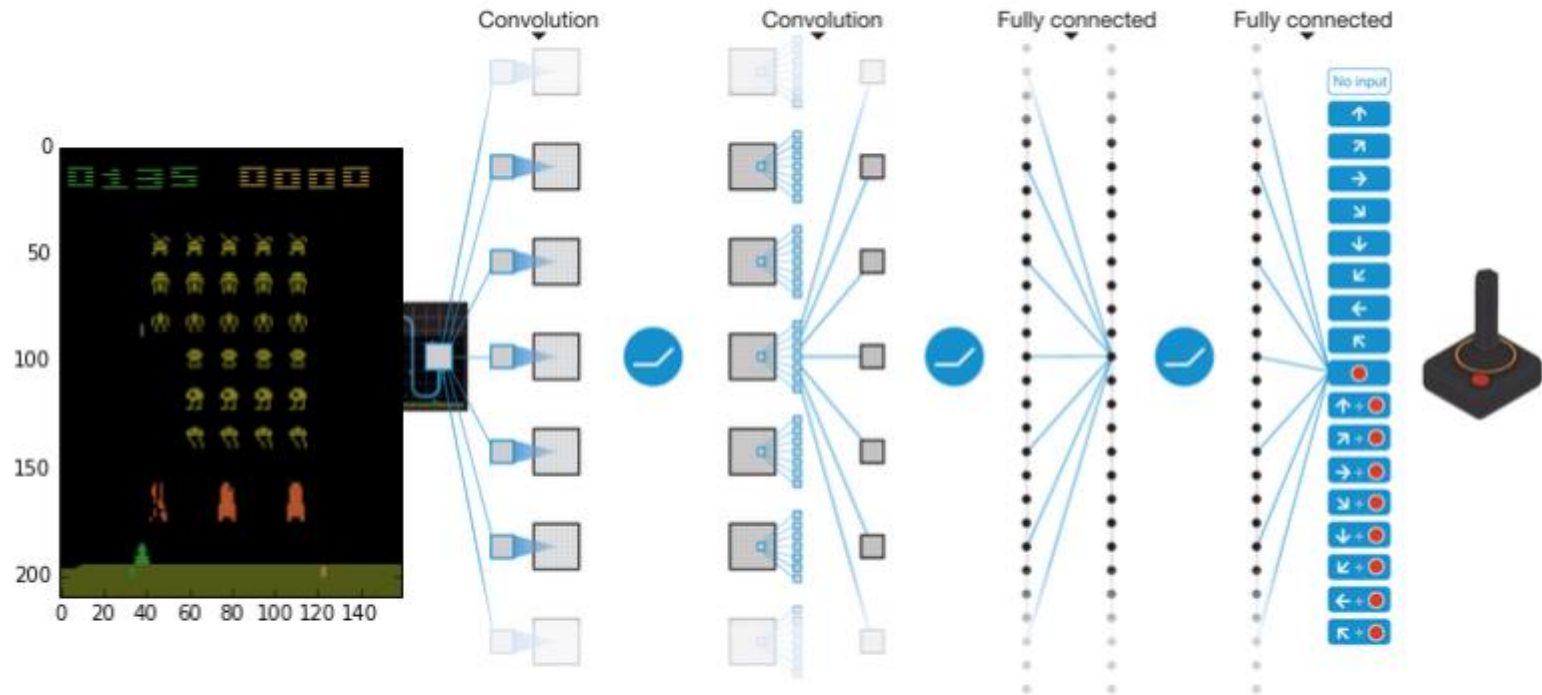


Given (\mathbf{s}, \mathbf{a})
Predict $Q(\mathbf{s}, \mathbf{a})$



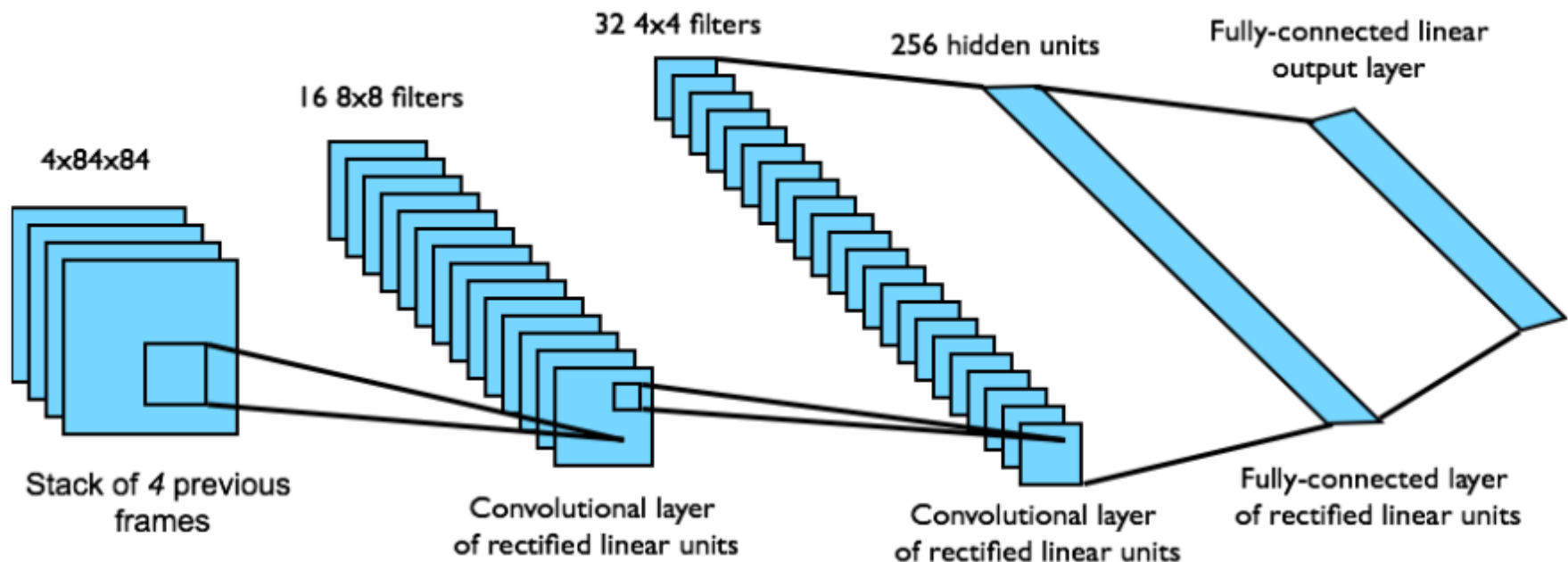
Given \mathbf{s} predict all q-values
 $Q(\mathbf{s}, \mathbf{a}_0)$, $Q(\mathbf{s}, \mathbf{a}_1)$, $Q(\mathbf{s}, \mathbf{a}_2)$

From theory to practice: DQN case

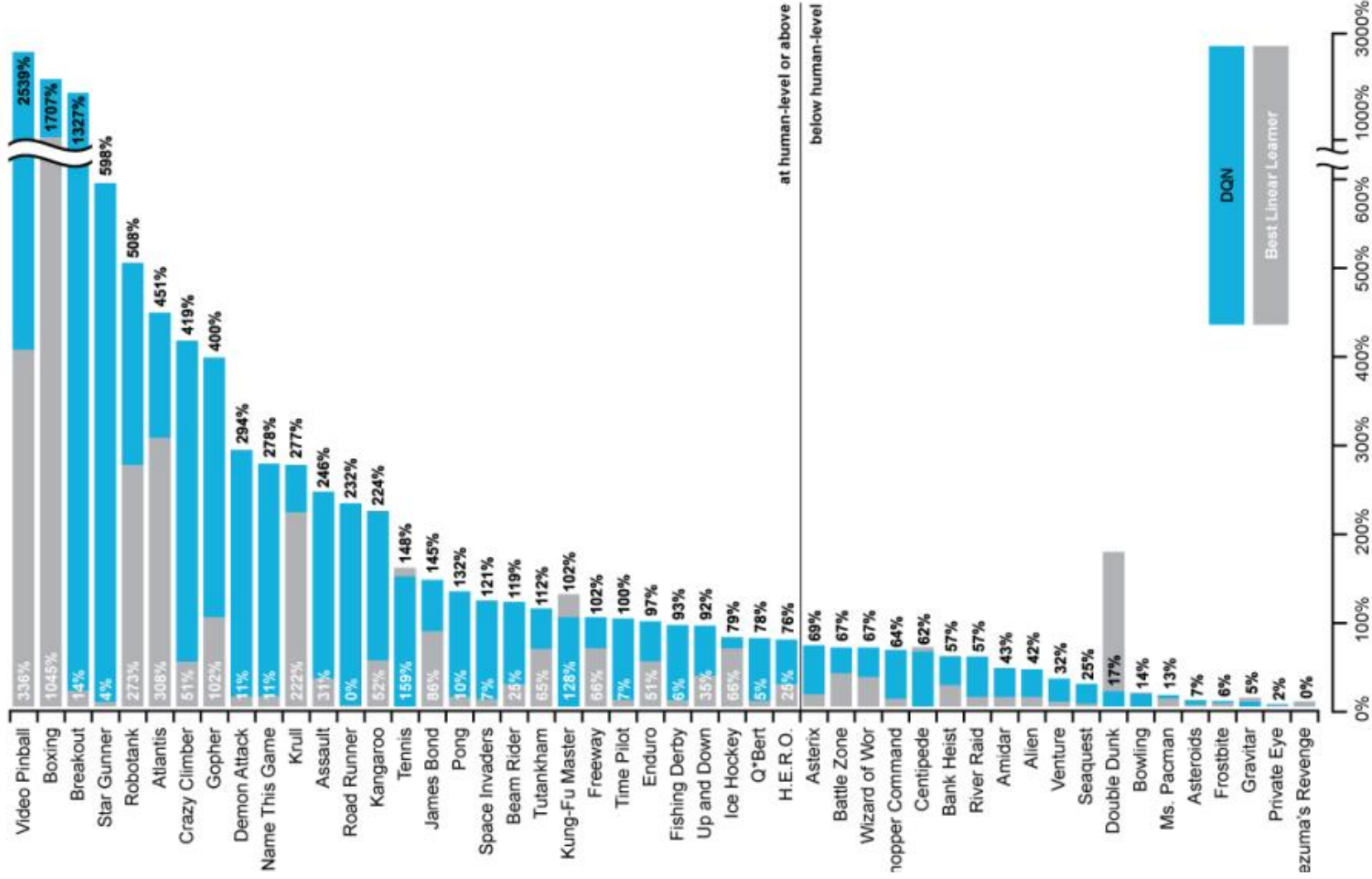


DQN: Atari

- End-to-end learning of values $Q(s, a)$ from pixels s
- Input state s is stack of raw pixels from last **4** frames
- Output is $Q(s, a)$ for **18** joystick/button positions
- Reward is change in score for that step



DQN results on Atari



DQN: under the hood

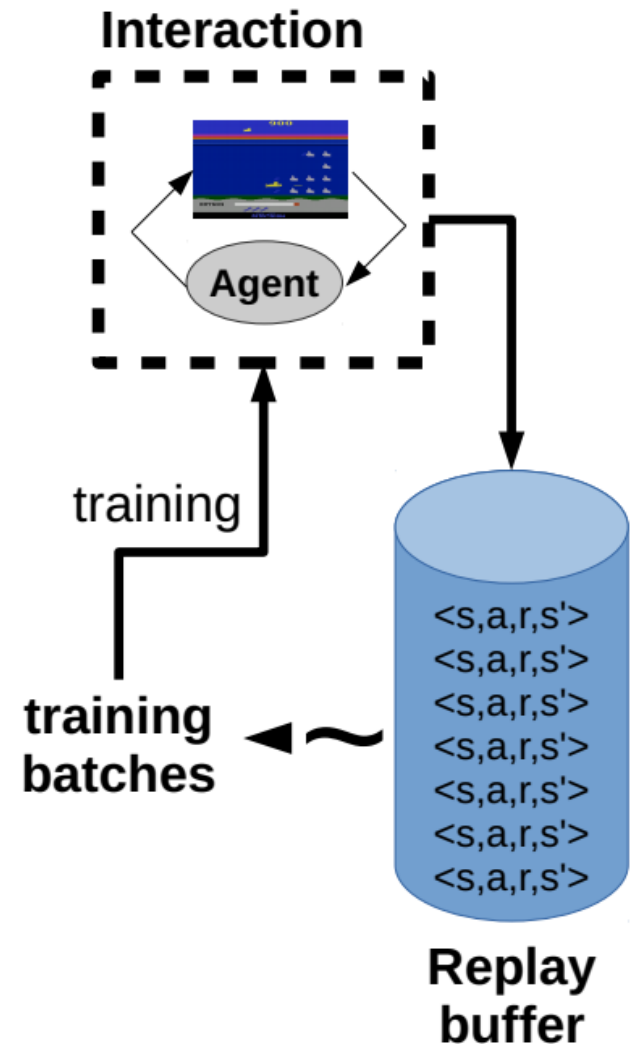
- DQN uses **experience replay** and **fixed Q-targets**:
 - Take action \mathbf{a}_t according to ϵ -greedy policy
 - Store transition $(\mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1})$ in replay memory \mathbf{D}
 - Sample random mini-batch of transitions $(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')$ from \mathbf{D}
 - Compute Q-learning targets w.r.t. old, fixed parameters \mathbf{w}^-
 - Optimize MSE between Q-network and Q-learning targets

$$\mathcal{L}_i(\mathbf{w}_i) = \mathbb{E}_{\mathbf{s}, \mathbf{a}, r, \mathbf{s}' \sim \mathcal{D}_i} \left[\left(r + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}'; \mathbf{w}_i^-) - Q(\mathbf{s}, \mathbf{a}; \mathbf{w}_i) \right)^2 \right]$$

- Using variant of stochastic gradient descent

Experience Replay

- **Idea:** store several past interactions $\langle s, a, r, s' \rangle$
- Train on random subsamples
- **Any +/- ?**



DQN: Atari Breakout

- <https://www.youtube.com/watch?v=TmPfTpjtdgg>