# CS 446/ECE 449: Machine Learning

A. G. Schwing

University of Illinois at Urbana-Champaign, 2020

L29: Q-learning

**Goals of this lecture**

- Extending MDPs
- Getting to know Q-learning

**Recap so far:** Known MDP

- To compute $V^*$, $Q^*$, $\pi^*$: use **policy/value iteration or exhaustive**
- To evaluate fixed policy $\pi$: use policy evaluation

But what if:

- Transition probabilities and rewards are not known
- No model available (model free RL)

Bellman optimality principle:

$$Q^*(s, a) = \sum_{s' \in \mathcal{S}} P(s' \mid s, a) \left[ R(s, a, s') + \max_{a' \in \mathcal{A}_{s'}} Q^*(s', a')) \right]$$

Idea:

- Run a simulator to collect experience tuples/samples $(s, a, r, s')$
- Approximate transition probability using samples

Algorithm sketch for Q-learning:

- Obtain a sample transition $(s, a, r, s')$
- Obtained sample suggests:

$$Q(s, a) \approx y_{(s,a,r,s')} = r + \max_{a' \in \mathcal{A}_{s'}} Q(s', a')$$

- To account for missing transition probability we keep running average

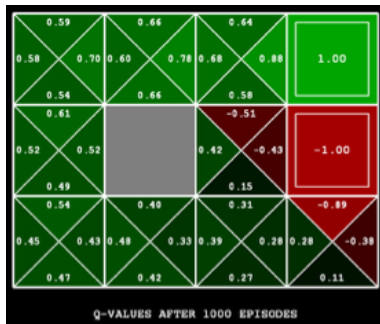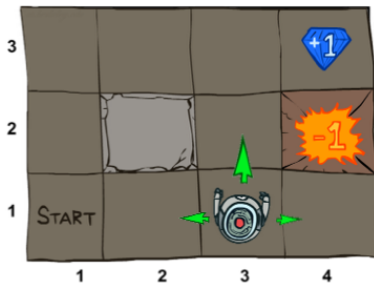$$Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha y_{(s,a,r,s')}$$

Summary:

- Known MDP
  - To compute $V^*$, $Q^*$, $\pi^*$: use value/policy iteration
  - To evaluate fixed policy $\pi$: use policy evaluation

- Unknown MDP: Model free
  - To compute $V^*$, $Q^*$, $\pi^*$: use Q-learning
  - To evaluate fixed policy $\pi$: use value learning

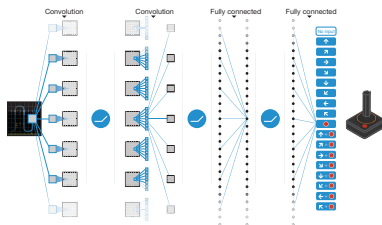Our current Q-learning works fine in case we have a tabular environment:

How to use this idea for playing Atari games?

- What are the actions?
- What are the rewards?
- What are the states?

How about $Q_\theta(s, a)$ as a neural net which takes images as input and produces a number for each of the possible actions?

Deep Q-learning algorithm (Deep Q-networks (DQN)):
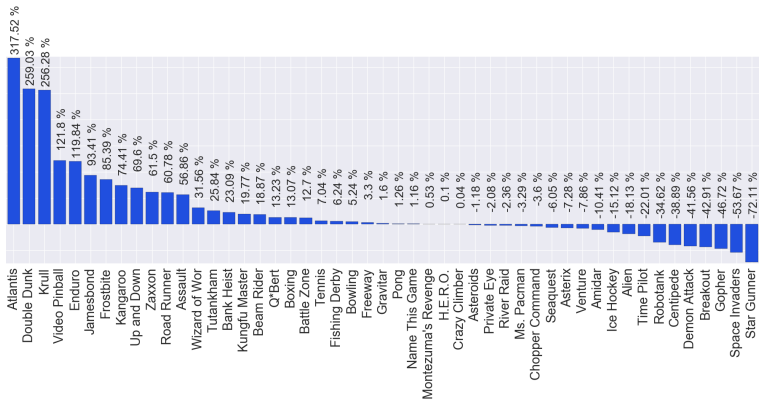
Given dataset $\mathcal{D} = \{(s_j, a_j, r_j, s_{j+1})\}$:

- Sample minibatch $\mathcal{B} \subseteq \mathcal{D}$
- Compute target $y_j = r_j + \gamma \max_a Q_{\theta^-}(s_{j+1}, a) \quad \forall j \in \mathcal{B}$
- Use stochastic (semi-)gradient descent to optimize w.r.t. parameters $\theta$

$$\min_\theta \sum_{(s_j, a_j, r_j, s_{j+1}) \in \mathcal{B}} \left( Q_\theta(s_j, a_j) - y_j \right)^2$$

- Perform $\epsilon$-greedy action and augment $\mathcal{D}$

Results:

**Quiz:**

- What differentiates RL from supervised learning?
- What is a MDP?
- What to do if no transition probabilities are available?

**Important topics of this lecture**

- Getting a feeling for reinforcement learning
- Understanding how to use reinforcement learning

**What's next:**

Policy Gradient Methods