

Informatics 1 Cognitive Science

Lecture 28: Reinforcement Learning Part 2

Matthias Hennig

School of Informatics
University of Edinburgh
mhennig@inf.ed.ac.uk

Credit Assignment Problem

Temporal Difference Learning

Q-Learning

RL through the ages

Credit Assignment Problem

Recalling the Rescorla-Wagner Rule

V_t : associative strength between CS and US at time t . Update rule:

$$\underbrace{V_t}_{\text{new}} = \underbrace{V_{t-1}}_{\text{old}} + \alpha \underbrace{(R - V_{t-1})}_{\delta}$$

- R : received reward
- $\alpha \in (0, 1)$: learning rate
- $\delta = R - V_{t-1}$: **prediction error**

Known as **δ -rule**: update \propto surprise.

Second-Order Conditioning

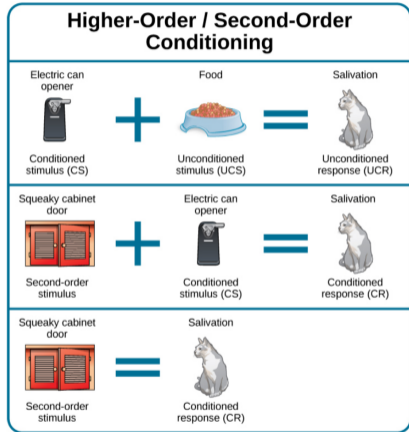


Image from Wikipedia.

1. A CS1 is first paired with reward.
 2. A second CS2 is then paired with CS1, without direct reward.
 3. CS2 acquires predictive value through CS1 and can trigger a conditioned response.
- Can the Rescorla-Wagner rule explain this?

Second-Order Conditioning

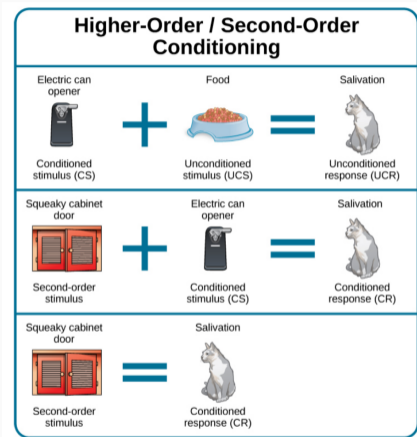


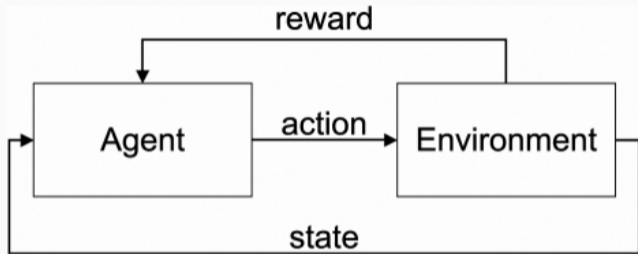
Image from Wikipedia.

The Rescorla-Wagner rule updates associations only when a US (reward) is present:

$$V_t = V_{t-1} + \alpha(R - V_{t-1})$$

- In second-order conditioning, CS2 is paired with CS1 **without any reward** ($R = 0$).
- With $R = 0$, the rule will **extinguish** any association rather than build one.
- It cannot transfer predictive value from CS1 to CS2.

Credit Assignment Problem



How to know which past action was responsible for an observed outcome? This is called the **temporal credit assignment problem**.

Temporal Difference Learning

Total Future Reward (in a finite trial)

Use the Rescorla-Wagner rule with stimulus $u(t)$ (0,1), reward $r(t)$ and expected reward $v(t)$ for trial time steps $t = 0, \dots, T$.

The total future reward from time t is the sum of all future rewards up to the end of the trial:

$$v(t) = \langle \sum_{\tau=t+1}^T r(\tau) \rangle$$

$\langle \rangle$ is the average over many trials. To compute this for a single stimulus, we use:

$$v(t) = \sum_{\tau=0}^t w(\tau) u(t - \tau)$$

$w(\tau)$ is the weight the stimulus from τ time steps ago contributes to the current reward prediction. How do we estimate $w(\tau)$?

Temporal Difference Rule

Goal: Learn the weights $w(\tau)$ so that the prediction

$$v(t) = \sum_{\tau=0}^t w(\tau)u(t - \tau)$$

matches the total future reward.

TD update rule:

$$w(\tau) \rightarrow w(\tau) + \alpha\delta(t)u(t - \tau)$$

$$\delta(t) = r(t) + v(t + 1) - v(t)$$

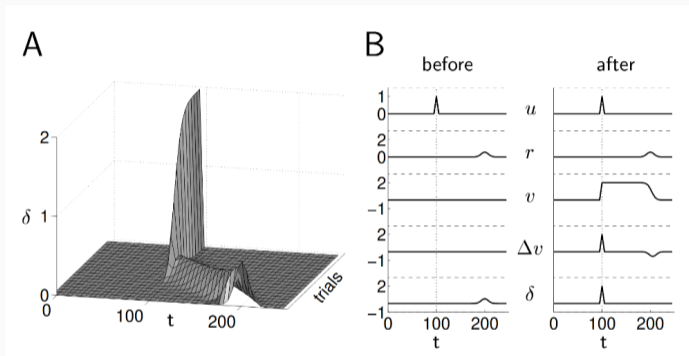
Temporal Difference Rule

The TD update rule:

$$\begin{aligned}w(\tau) &\rightarrow w(\tau) + \alpha\delta(t)u(t - \tau) \\ \delta(t) &= r(t) + v(t + 1) - v(t)\end{aligned}$$

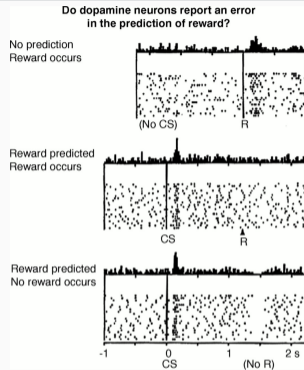
- The weights are adjusted proportional to the TD error $\delta(t)$.
- $v(t + 1) - v(t)$ compares two successive reward predictions.
- High reward / higher prediction in future: $\delta(t) > 0$
- Low reward / lower prediction in future: $\delta(t) < 0$

Temporal Difference Learning



Temporal Difference learning shifts the prediction-error signal earlier in time.

Reward Learning in the VTA



(1) neurons signal a prediction error. (2) Once learned, neurons signal reward at the time of the stimulus (Schultz 1998, Predictive reward signal of dopamine neurons. Journal of Neurophysiology.)

From stimuli to states

TD learning predicts reward from the recent stimulus history:

$$v(t) = \sum_{\tau=0}^t w(\tau) u(t - \tau)$$

Now we introduce a **state-based** description: if the agent is in state s_t , we learn a state value $V(s_t)$.

Temporal-difference error:

$$\begin{aligned}\delta_t &= r_{t+1} + \gamma V(s_{t+1}) - V(s_t) \\ V(s_t) &\leftarrow V(s_t) + \alpha \delta_t\end{aligned}$$

$\gamma \in (0, 1]$ discounts future rewards; α is the learning rate.

Q-Learning

From states to actions: Bellman Equation and Q-values

The state values now allow learning which **actions** – moving between states – to take in each state to maximise reward.

Bellman equation:

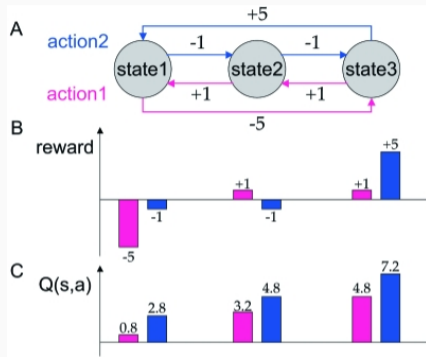
$$V_{k+1}(s) \leftarrow \max_a (R + \gamma V_k(s'))$$

s : current state, s' : next state, $\gamma \in (0, 1]$: discount factor, chocolate now or chocolate tomorrow? This finds the **optimal policy** by maximizing expected future reward.

Q-learning has the same aim, but is more efficient: collect the expected reward of **state-action pairs**:

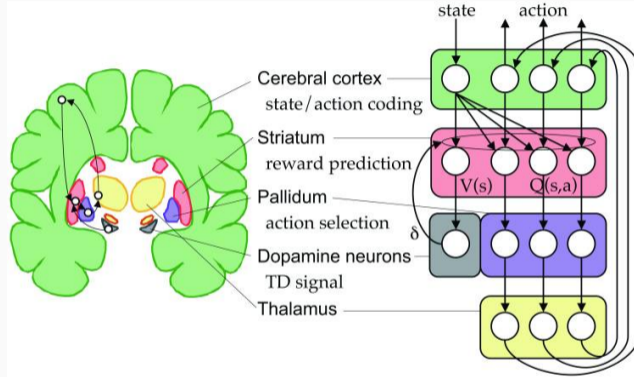
$Q(s, a) \approx$ expected future reward if we take action a in state s

Q-Learning: TD Learning + choosing actions



Discount factor: $\gamma = 0.8$. For state 2, it is now better to accept the negative reward first, as more reward is on the horizon later.

Action Learning in the Brain



Dopamine neurons in the midbrain: prediction errors δ ; Striatum: Value; Cortex: World

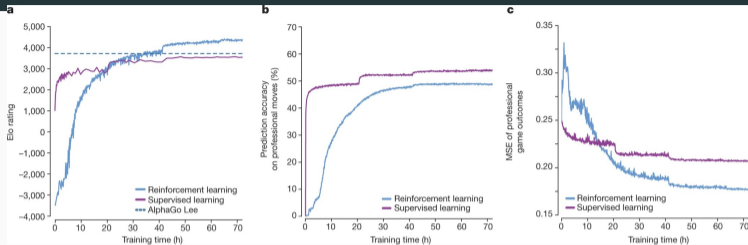
RL through the ages

Reinforcement Learning Successes



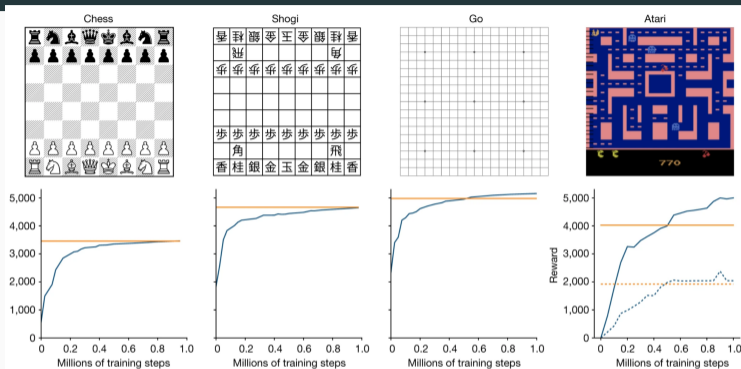
- RL from pixels: human-level Atari video game playing.
- AlphaGo defeated Lee Sedol (2016); successor AlphaZero reached superhuman play via self-play.
- RLHF / preference optimization fine-tunes large language models.

AlphaGo Zero (2017)

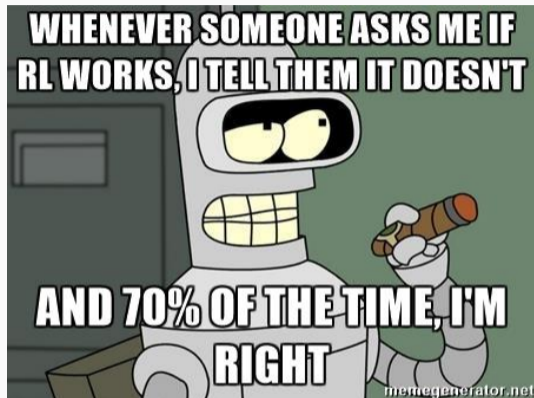


- Knows game rules, but trained by playing against itself (64 GPU workers and 19 CPU parameter servers, \$25M).
- Best Models are used as new opponents for self-play.
- 0.4s thinking time/move, 4.9 million games played; 216,000 moves/day. About 3 days training in total.
- Comparison to human-trained supervised model (move prediction).

MuZero (2020)



- Does not know game rules.
- Uses an internal (hidden, latent) model.
- Trained to predict actions, values and rewards.
- Board games: 16 TPUs training / 1000 TPUs self-play. 5 hours training.
- Comparison to AlphaZero and human-level (Atari) performance.



"Deep Reinforcement Learning Doesn't Work Yet" (Irpan 2018)

<https://www.alexirpan.com/2018/02/14/rl-hard.html>

- TD learning generalizes the δ -rule to learn predictions over time.
- Q-learning extends TD learning to choose actions via state–action values.
- Deep RL achieves impressive performance in learning from actions (but is far from human-level learning).