# Autonomous Learning

Installing your programming environment

Simon Reichhuber

May 27, 2020

University of Kiel, Summer Term 2020

# TABLE OF CONTENTS

# Submissions

- Reinforcement Learning



- StarCraft II API (Blizzard & Google)

- Upload your solution to the link of your Group ($\leftarrow$ posted in the forum)
- Submit a zip-file containing the code and an explainatory video: *<Group name>_Assignement<Assignement no.>.zip*

Recommendation: use Linux

- Besides, some specific problems visualisations cannot be switched off on Windows/Mac
- Usually packages are easier to install on Linux

Remarks to coding style

- It is not forbidden to use object-oriented programming.
- Use folders to increase readability (e.g. assignement1, task1, etc.)
- Comment important code
- Use telling variable names, like *update_state_values()* (not *do_calcuation()*, or *get_result()*, etc. )

Important: Use debugging tools of your IDE

- DeepMind's Python Framework
- Reinforcement Learning Environment
  - receive Observation
  - take Action
- github: www.github.com/deepmind/pysc2

- Knowledge limited to the scope of human players
- No chance against simple KIs
- For hard coding KIs too complex

Stand 2019:
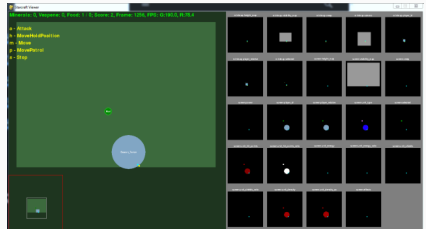AlphaStar wins 10:1 against professional StarCraft 2 players

- Isolated view of certain elements of the game
- Focused scenarios on small maps
- Test/Learn actions and/or game mechanics

**Minigames**

- MoveToBeacon
- CollectMineralShards
- FindAndDefeatZerglings
- DefeatRoaches
- DefeatZerglingsAndBanelings
- CollectMineralsAndGas
- BuildMarines

Apply different Reinforcement Learning algorithms on the minimap MoveToBeacon

Marine:

```
1 def get_marine(obs):
2     marine = next(unit for unit in obs.observation.
          feature_units
3                   if unit.alliance == features.
                       PlayerRelative.SELF)
4     return marine
```

Beacon:

```
1 def get_beacon(obs):
2     beacon = next(unit for unit in obs.observation.
          feature_units
3                   if unit.alliance == features.
                       PlayerRelative.NEUTRAL)
4     return beacon
```

Position:

```python
1  def getCoord(unit):
2      x = unit.x
3      y = unit.y
4      return x, y
```

Keine Aktion:

```
1  return actions.FUNCTIONS.no_op()
```

Select Unit:

```
1  return actions.FUNCTIONS.select_army("select")
```
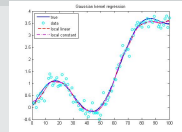
Move Selected Units:

```
1  return actions.FUNCTIONS.Move_screen("now", (x, y))
```
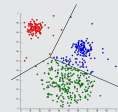
# Reinforcement Learning

## Supervised Learning

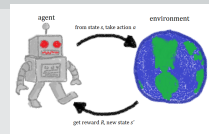$\rightarrow$ problem-driven (Regression / Classifiction)



## Unsupervised Learning

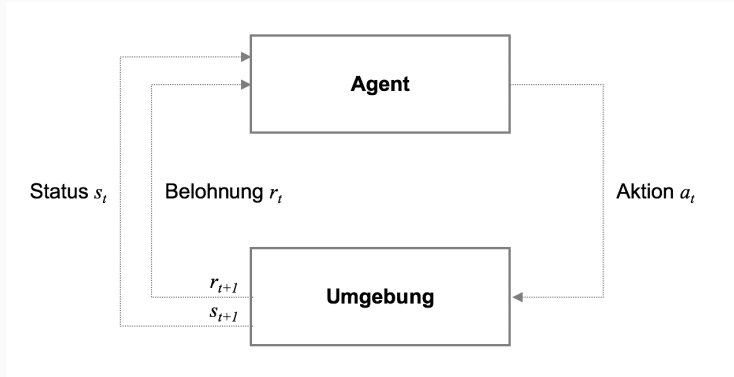$\rightarrow$ data-driven (Clustering)



## Reinforcement Learning

$\rightarrow$ react to environment

**Reinforcement Learning**

"Reinforcement learning is learning what to do—how to map situations to actions—so as to maximize a numerical reward signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate reward but also the next situation and, through that, all subsequent rewards. These two characteristics—trial-and-error search and delayed reward—are the two most important distinguishing features of reinforcement learning."

*Richard S. Sutton and Andrew G. Barto*

1. Agent
2. Environment
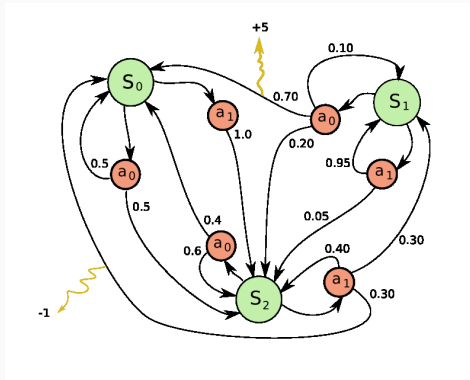3. State
4. Action
5. Reward

**Goal**

Find relation of action and expected reward and maximise rewards by taking several actions.

Problem: Follow known strategies (exploitation) or explore new unkown strategies which might result in a bigger reward?

$\Rightarrow$ Exploration-Exploitation Dilemma

## Markov Property

The future is independent of the past if you know the presence.

- 4 actions
- Each action reults in a reward of -2
- Coloured states are final states with a specific reward (reward 100/-100)



- With probability 0.8 the agent moves to the intended state; with probability 0.1 it takes one of the other two orthogonal actions.

## Notation

- $S = s_1, \ldots, s_n$ is the set of all states
- $A = a_1, \ldots, a_n$ is the set of all actions
- $T : S \times A \times S \to [0, 1]$ is the transition function. $T(s, a, s')$ represents the probability that an agent in state $s$ moves to $s'$ by taking action $a$.
- $R : S \times A \times S \to \mathbb{R}$ is the reward function. $R(s, a, s')$ assigns the transition from state $s$ to state $s'$ by taking action $a$ a reward.

**Notation**

- $\pi(s)$ is a policy that maps the state $s$ to an optimal action.

- $\gamma \in [0, 1]$ is the discount factor. All possibly future rewards are multiplied after each action with $\gamma$. By this, the agent shall prefer short paths. With $\gamma = 1$ the length of the path is arbitrary if the reward is equal along the paths. With $\gamma = 0$ only the reward of the next step is taken into account.

**V-Value**

$V^*(s)$= Estimated total reward for an agent die der Agent starting at a specific state $s$ and taking the optimal actions afterwards.
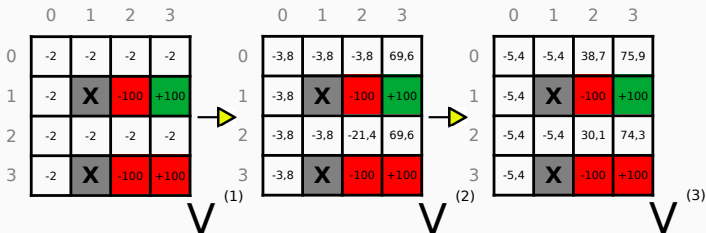
$$V^*(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s')[R(s, a, s') + \gamma * V^*(s')]$$

**Meaning**

What is the value of a state $s$?

- Temporally limited MDP
- $V^{(k)}$ = Total reward, only *k* steps left

$$V^{(k)}(s) = \max_{a \in A} \sum_{s' \in S} T(s, a, s')[R(s, a, s') + \gamma * V^{k-1}(s')]$$

Aktion: South

$$0.8 \times \left(-2 + 0.9 \times 100\right) = 70.4$$

$$0.1 \times \left(-2 + 0.9 \times -2\right) = -0.38$$
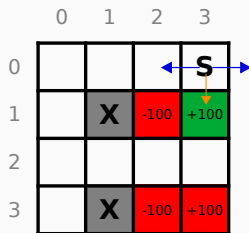
$$0.1 \times \left(-2 + 0.9 \times -2\right) = -0.38$$

$$\sum_{s' \in S} \left[70.4; -3.04; -3.04\right] = 69.6$$



Aktion: West $= 5,38$; Aktion: North $= -3,8$; Aktion: East $= 5,38$

$$\max_{a \in A}\left[69,6; 5,38; -3,8; 5,38\right] = 69,6$$

$$V^{(2)}(0,3) = 69,6$$

- max Reward of the MDP nach bounded above
- $V^k$ in the limit $k \to \infty$ converges
- For sufficiently large $k$ it holds $V^{(k)} = V^* \pm \epsilon$

- From the v-values $V^*$ one can derive a optimal policy $\pi^*$
- For each state take the action that maximises the reward

$$\pi^*(s) = \underset{a \in A}{\operatorname{argmax}} \sum_{s' \in S} T(s, a, s')[R(s, a, s') + \gamma * V^*(s')]$$

# Q-Learning

- Value iteration and policy extraction requires knowledge about $T(s, a, s')$ and $R(s, a, s')$
- Agent has to follow its own strategy from state to state
- Externally, the agent receives a reward from the environment

- Initially, randomly chosen V values $\rightarrow$ iterative improvement
- After every step, the agent analyses the reward and the possible future rewards and adapts its V value estimation approximatively.

$$V(s) \leftarrow V(s) + \alpha[r + \gamma V(s') - V(s)]$$

**Q-Values**

$Q(s, a)$ is the quality of action $a$ in state $s$

$$Q^*(s, a) = \sum_{s' \in S} T(s, a, s')[R(s, a, s') + \gamma V^*(s')]$$

explicit actualisation of the policy can be unnecessary if it is immediately learned implicitely (on-policy):

$$V(s) = \max_{a \in A} Q(s, a)$$

$$\pi(s) = \underset{a \in A}{\mathrm{argmax}} \, Q(s, a)$$

**Q-Learning**

$=$ Temporal Difference Learning auf Q-Values

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a' \in A} Q(s', a') - Q(s, a)]$$

- Q function as table
- For each Q value a single entry
- Initialise the values (randomly, with max reward, or zero)

| States \ Actions | NORTH | EAST | SOUTH | WEST |
|---|---|---|---|---|
| (3,0) | 31 | 26,4 | 25,7 | 26,4 |
| (2,0) | 37,6 | 28,1 | 26,5 | 32,3 |
| (2,1) | 28,1 | 32,2 | 28,1 | 31,2 |
| (2,2) | -69,9 | 39,6 | -69,9 | 1,5 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

- Exploration strategy
- Hyperparameter $0 \leq \epsilon \leq 1$
- On every step, the agent takes a random action with probability $\epsilon$
- $\epsilon$ decreases over training time

Initialize $Q(s, a)$ arbitrarily;

**repeat**

    Initialize $s$;

    **repeat**

        Chose $a$ from $s$ using policy derived from $Q$

         ($\epsilon$-greedy);

        Take action $a$, observe $r$, $s'$;

        $Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a \in A} Q(s', a) - Q(s, a)]$;

        $s \leftarrow s'$

    **until** $s$ *is terminal*;

**until** $Q$ *is converged*;

**Algorithm 1:** Q-Learning

# Bibliography

1. https://medium.freecodecamp.org/diving-deeper-into-reinforcement-learning-with-q-learning-c18d0db58efe
2. https://github.com/deepmind/pysc2
3. https://www.mathworks.com/matlabcentral/mlc-downloads/downloads/submissions/19564/versions/1/screenshot
4. https://2s7gjr373w3x22jf92z99mgm5w-wpengine.netdna-ssl.com/wp-content/uploads/2018/02/reinforcement-learning.png
5. https://towardsdatascience.com/reinforcement-learning-demystified-markov-decision-processes-part-1-bf00dda41690
6. https://hciweb.iwr.uni-heidelberg.de/system/files/private/downloads/541645681/damma-learning-report.pdf